# Multiscale recognition of legume varieties based on leaf venation images

Mónica G. Larese [a,b,*], Ariel E. Bayá [a], Roque M. Craviotto [b], Miriam R. Arango [b], Carina Gallo [b], Pablo M. Granitto [a]

[a] CIFASIS, French Argentine International Center for Information and Systems Sciences, UAM (France)/UNR-CONICET (Argentina), Bv. 27 de Febrero 210 Bis, 2000 Rosario, Argentina
[b] Estación Experimental Oliveros, Instituto Nacional de Tecnología Agropecuaria, Ruta Nacional 11 km 353, 2206 Oliveros, Santa Fe, Argentina

## ARTICLE INFO

## ABSTRACT

In this work we propose an automatic low cost procedure aimed at classifying legume species and varieties based exclusively on the characterization and analysis of the leaf venation network. The identification of leaf venation patterns which are characteristic for each species or variety is not an easy task since in some situations (specially for cultivars from the same species) the vein differences are visually indistinguishable for humans. The proposed procedure takes as input leaf images acquired using a standard scanner, processes the images in order to segment the veins at different scales, and measures different traits on them. We use these features in combination with modern automatic classifiers and feature selection techniques in order to perform recognition. The process was initially applied to recognize three different legumes in order to evaluate the improvements over previous works in the literature, and then it was employed to distinguish three diverse soybean cultivars. The results show the improvements achieved by the usage of the multiscale features. The cultivar recognition is a more challenging problem, since the experts cannot distinguish evident differences in plain sight. However, we achieve acceptable classification results. We also analyze the feature relevance and identify, for each classifier, a small set of distinctive traits to differentiate the species and varieties.

© 2014 Elsevier Ltd. All rights reserved.

## 1. Introduction

Many works in the current literature deal with the problem of automatically identifying plants by means of foliar image analysis. One of the most common approaches consists in performing shape analysis of the leaves (Agarwal et al., 2006; Camargo Neto, Meyer, Jones, & Samal, 2006; Chaki & Parekh, 2012; Du, Wang, & Zhang, 2007; Im, Nishida, & Kunii, 1998; Solé-Casals, Travieso, Alonso, & Ferrer, 2008). Leaf color and texture can also be taken into consideration. In the work by Pydipati, Burks, and Lee (2006), color texture features of the leaves are used in combination with discriminant analysis to detect citrus diseases. Also, a combination of shape, texture and color features are used in the papers by Golzarian and Frick (2011) and Bama, Valli, Raju, and Kumar (2011).

However, in some practical situations there are not evident differences in the shape, size, color or texture features of the leaves for the plants under study. This is the case, for example, of plants that belong to several cultivars from the same species.

Since there exists correlation between leaf venation characteristics and leaf properties (such as damage and drought tolerance, among others) (Sack, Dietrich, Streeter, Sanchez-Gomez, & Holbrook, 2008; Scoffoni, Rawls, McKown, Cochard, & Sack, 2011), some works in the recent literature highlight the importance of analyzing the structure of the venation system as a means to perform leaf-based plant identification. In the paper by Park, Hwang, and Nam (2008), a content-based image retrieval system is proposed which analyzes the venation of a leaf sketch drawn by the user as an initial categorization, and then uses shape features to find similar leaves existing in the database. On the other hand, Clarke et al. (2006) and Valliammal and Geethalakshmi (2011) propose new methods for leaf vein segmentation. However, neither work includes any characterization or recognition tasks. Recently, Du, Zhai, and Wang (2013) propose a method based on fractal dimension features computed both on the veins and the leaf outline, and employ k-nearest neighbors to perform classification on different leaves. However, these leaves are visually very different and belong to very different families. Additionally, the computed

* Corresponding author at: CIFASIS, French Argentine International Center for Information and Systems Sciences, UAM (France)/UNR-CONICET (Argentina), Bv. 27 de Febrero 210 Bis, 2000 Rosario, Argentina. Tel.: +54 (0) 341 4237248x325; fax: +54 (0) 341 4237248x301.

E-mail addresses: larese@cifasis-conicet.gov.ar (M.G. Larese), baya@cifasis-conicet.gov.ar (A.E. Bayá), rcraviotto@correo.inta.gov.ar (R.M. Craviotto), marango@correo.inta.gov.ar (M.R. Arango), cgallo@correo.inta.gov.ar (C. Gallo), granitto@cifasis-conicet.gov.ar (P.M. Granitto).

features do not provide the experts with a simple vein description and may not lead to a human direct interpretation.

In a recent work, Price, Symonova, Mileyko, Hilley, and Weitz (2011) developed an interactive graphic tool named LEAF GUI, aimed at thresholding, cleaning and segmenting stained leaf vein and areole images in a user-assisted way. In addition, the software allows to extract several measures which are automatically computed on these structures. The segmentation algorithms require that the visibility of the veins were previously enhanced by means of X-ray techniques, chemical or biological clearing, or back-lit scanning. LEAF GUI does not include any feature selection algorithm or any plant classification/recognition procedure.

Agricultural specialists require, in many situations, to identify which species a certain batch of plants corresponds to. Depending on the species and growth stage of the plants, leaves, flowers, fruits and/or seeds can be used in conjunction to recognize the species. But in many other situations, this is not possible. For example, if the goal is to differentiate diverse cultivars/varieties from the same species, all the previous mentioned characteristics may be visually the same. One possibility is to perform DNA analysis to accurately determine the variety, but this method is expensive. On the other hand, we propose to investigate the possibility of searching for distinctive venation patterns that could uniquely identify the varieties using a low cost procedure based on an image analysis and machine learning system.

We showed in a previous work (Larese et al., 2014) that it is possible to recognize different species using exclusively information from the leaf veins. The motivation of the present work is to extend the analysis to the more difficult problem of recognizing varieties from the same species. We search for the existence of distinctive leaf vein patterns for different cultivars, when all the other leaf characteristics (e.g., shape, color and texture) are similar. If the plants under study have different physiological characteristics (e.g., drought tolerance), there is a chance that these properties can be reflected in their veins even if the leaves look similar. In this work we propose an automatic low cost procedure aimed at segmenting and characterizing the leaf veins of plants from the same family. An automatic procedure is desirable since it provides reliability, reproducibility and economy, besides of providing a solution to a problem which is not easily solved by the human experts, as it is the cultivars recognition.

Since this problem is more difficult than separating different species, we propose to measure vein traits from images at different scales. We first try the procedure on the simpler problem of species recognition, showing that the new approach improves the results reported in our previous work (Larese et al., 2014). Next, we analyze the cultivar recognition problem.

We use three legume species and three soybean cultivars in order to perform the species and variety recognition, respectively. The leaves are acquired using a standard flatbed scanner. We perform the automatic plant recognition by means of measuring and classifying morphological traits from central patches extracted from the previously segmented venation system, i.e., no leaf shape, color or texture information is considered. We also analyze the distinctive vein characteristics for each class.

For this purpose, we start by performing segmentation using the Unconstrained Hit-or-Miss Transform (UHMT) and adaptive image thresholding in order to extract the veins at several image scale levels. The UHMT is a mathematical morphology operator useful to perform template matching. It extracts all the pixels which follow a certain foreground and background neighboring configuration.

After segmentation, we compute several morphological measures on the segmented veins at the different scales, and use them as features in the classification process. The recognition is performed resorting to three different classifiers, namely, Random Forests (Breiman, 2001), Support Vector Machines with Gaussian kernel (Vapnik, 1995) and Penalized Discriminant Analysis (Hastie, Buja, & Tibshirani, 1995). Recursive Feature Elimination (Guyon, Weston, Barnhill, & Vapnik, 2002) is also used in combination with the three classifiers in order to estimate the importance of the input variables in the classification process for the different species and varieties.

The analysis is performed on two different problems. First of all, we consider the discrimination between three classes of legumes, namely soybean (*Glycine max (L) Merr*), red and white beans (*Phaseolus vulgaris*). Red and white beans have very similar leaves, which are slightly darker for the former. However, in this work we do not consider color information, but only morphological features of the veins obtained from the gray scale images.

The second problem consists of identifying three different cultivars of soybean. This task is more challenging by far, since the differences in the veins are not obvious to the human experts. Automatic classification would come to solve this issue in an inexpensive way. Additionally, the procedure would highlight relevant distinctive vein features for each cultivar, and possibly help to relate these differences to variety adaptation.

The rest of the paper is organized as follows. In Section 2.1 we describe the leaf images dataset. Sections 2.2 and 2.3 summarize the segmentation procedure that we employed to extract the leaf venation system. We detail the measures computed on the segmented veins in Section 2.4. We briefly describe the classification and feature selection algorithms in Section 2.5. We present and discuss the results in Section 3, where we assess the performance of the procedure and analyze the relevant features. Finally, we draw some conclusions in Section 4.

## 2. Materials and methods

### 2.1. Leaf images dataset

The dataset used in this paper is composed by a total number of 866 color leaf images provided by Instituto Nacional de Tecnología Agropecuaria (INTA, Oliveros, Argentina). The dataset is divided in the following way: 422 images correspond to soybean leaves (198 belong to cultivar 1, 176 belong to cultivar 2, and 48 belong to cultivar 3), 272 images are from red bean leaves and 172 from white bean leaves. They are the images of the first foliage leaves (preformed in the seed) of 433 specimens (211 soybean plants, 136 red bean plants and 86 white bean plants). First foliage leaves were selected for the analysis, after 12 days of seedling grow, since their characteristics are less influenced by the environment. We did not use any chemical or biological procedure to physically enhance the leaf veins. Instead, a fast, inexpensive and simple imaging procedure was used: the leaves were acquired using a standard flatbed scanner (Hewlett Packard Scanjet-G 3110) at a resolution of 200 pixels per inch, and the images were stored as 24-bit uncompressed TIFF images. We scanned the abaxial surfaces of the leaves, since veins appear stronger on this side and can be considerably better observed.

### 2.2. Unconstrained Hit or Miss Transform (UHMT)

The UHMT is an extension of the Hit-or-Miss Transform (HMT) for gray scale images (Soille, 1999). It extracts all the pixels matching a certain foreground and background neighboring configuration. A composite structuring element **B** is employed, which is a disjoint set formed by one structuring element that specifies the foreground configuration, $B_{fg}$, and one structuring element for the background setting, $B_{bg}$. The origin of the composite structuring element matches the foreground.

The UHMT is defined as

$$UHMT_{\mathbf{B}}(Y)(y) = \max\left\{\varepsilon_{B_{fg}}(Y)(y) - \delta_{B_{bg}}(Y)(y), 0\right\}, \qquad (1)$$

where $Y$ is a gray scale image with set of pixels $y$ and $\mathbf{B}$ is a composite structuring element. It can be computed as the difference between an erosion with $B_{fg}$, $\varepsilon_{B_{fg}}(Y)(y)$, and a dilation with $B_{bg}$, $\delta_{B_{bg}}(Y)(y)$, if $\delta_{B_{bg}}(Y)(y) < \varepsilon_{B_{fg}}(Y)(y)$. Otherwise it equals 0.

## 2.3. Vein segmentation

The vein segmentation procedure is shown in Fig. 1 for a white bean leaf. In the following, we describe each image processing step in detail.

Since we are only interested in the vein morphological patterns, we removed from the images all the color information by converting the RGB images to gray scale. In order to perform this, we followed a standard procedure. We calculated the luminance component ($Y$) as a weighted sum of the three color channels ($R, G, B$), i.e., $Y = 0.299R + 0.587G + 0.114B$ (Umbaugh, 2005).

Then, we thresholded the gray scale image $Y$ by means of the automatic iterative threshold selection algorithm (Sonka, Hlavac, & Boyle, 2008) and filled its holes using morphological reconstruction (Soille, 1999). After deleting all the connected components except the largest one, we got a binary mask for the leaf.

We computed the UHMT on 5 different sized versions of $Y$, namely at 100%, 90%, 80%, 70% and 60%. Next we resized back each resulting UHMT to the original image size. We summed these five resized UHMTs to obtain the combined UHMT, which highlights both small and large visible veins simultaneously.
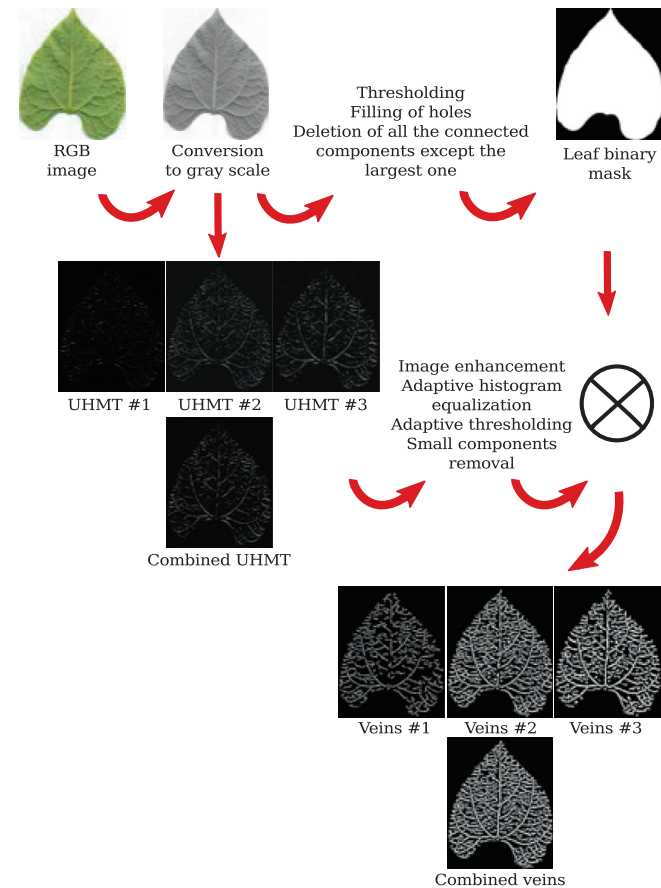


**Fig. 1.** Vein segmentation procedure.

On the other hand, we also preserved the resized UHMTs calculated at 100%, 80% and 60%, namely UHMT #1, UHMT #2 and UHMT #3 (i.e., UHMT #1 is the UHMT computed on the original gray scale image $Y$). Each UHMT is intended to highlight a different level of vein detail.

We used four composite structuring elements (foreground and background configurations) to detect the leaf veins in four directions (vertical, horizontal, +45 and −45 degrees). They are depicted in Fig. 2. Each composite structuring element describes foreground and background searched hits (in red and green, respectively).

Then, we enhanced the contrast of the four obtained UHMTs (the combined UHMT and UHMTs #1, #2 and #3). We performed adaptive histogram equalization and adaptive thresholding, and removed all the connected components with less than 20 pixels. Finally, we multiplied the four resulting UHMTs by the previously computed leaf binary mask to obtain four approximations of the veins (combined veins + 3 scales, i.e., veins #1, veins #2 and veins #3).

## 2.4. Vein measurements

After segmenting the venation system, we measured several traits on the veins and the areoles. Since our goal is to perform classification considering only the veins morphology, we avoided the influence of the leaf shape by cropping a centered $100 \times 100$-pixel patch for each one of the 4 vein images (the three scales and the combined veins), and measured all the features on these patches.

In this work we adapted LEAF GUI (Price et al., 2011) measures to extract a set of features of interest for veins and areoles. For our particular problem aimed at leaf classification, individual vein/areole measures are not suitable. For this reason, we computed the median, minimum and maximum feature values for veins and areoles where it was appropriate. We measured the 52 traits described in Table 1.

These features were computed for each one of the three scales and the combined veins patches. Altogether, they become a feature vector of 208 components (52 features×4 patches) per leaf image.

All the features we computed are rotational independent when considering the whole leaf, except for the edge orientation measures, namely VmO, VMO and VMeO. These features measure the angle between the $x$-axis of the image and the major axis of the ellipse having the same second moments as the vein, so if the leaf is rotated, these 3 features change.

However, since we measured the features on a square patch at the center of the leaf, all the measures would be affected if the leaf is rotated. This could be avoided by taking a circular patch instead of a square one, although the rotational dependence remains for the 3 orientation features mentioned above.

In this paper, we dealt with this issue by scanning all the leaves in the same vertical position, thus avoiding significant rotation influences.

For future general applications, a preprocessing module can be added in which all the images are previously corrected for rotation taking as a reference the longitudinal primary vein.

## 2.5. Classification algorithms

We considered 3 different classifiers, namely Random Forests, Support Vector Machines with Gaussian kernel and Penalized Discriminant Analysis. Each one of them is briefly described in the following subsections.

### 2.5.1. Random Forests (RF)

Random Forests (RF) (Breiman, 2001) is a state-of-the-art ensemble algorithm where the individual classifiers are a set of de-correlated trees. They perform comparably well to other
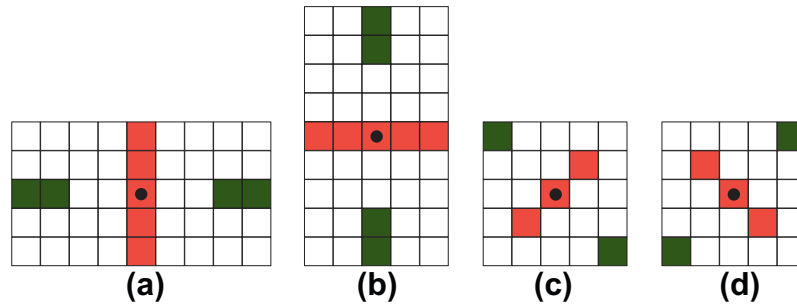
**Fig. 2.** The four pairs of flat composite structuring elements used in the UHMT computation to detect veins in four directions: (a) vertical, (b) horizontal, (c) +45 degrees, and (d) −45 degrees. Foreground pixel configurations are depicted in red while background pixel configurations are in green. The center of the composite structuring element is marked with a black dot. (For interpretation of the references to colour in this figure legend, the reader is referred to the web version of this article.)

**Table 1**
Set of features measured on the leaf veins after segmentation.

| Feature # | Code | Description |
|---|---|---|
| 1 | VNE | Total number of edges, i.e., estimated veins. |
| 2 | VNN | Total number of nodes. The number of connecting nodes between edges. |
| 3 | VTNL | Total network length. Total distance (in mm) along the skeleton of the vein image patch. |
| 4/5/6 | VMeL/VmL/VML | Median/min/max edge length. The edge length (in mm) is the distance along the skeleton of a vein. |
| 7/8/9 | VMeW/VmW/ VMW | Median/min/max edge width. The edge width (in mm) is the mean of the doubled distances between each skeleton pixel of the current edge and the nearest non-vein pixel, i.e., areole pixel. |
| 10/11/12 | VMeA/VmA/VMA | Median/min/max edge 2D area. The edge 2D area (in mm$^2$) is the sum of the widths at every skeleton pixel of the current edge times the length of one pixel. |
| 13/14/15 | VMeSA/VmSA/ VMSA | Median/min/max edge surface area. The surface area (in mm$^2$) of the cylinder centered at the edge skeleton is computed as the sum of the individual surface areas for each skeleton pixel of the current edge, as $\sum_i SA_i = 2\pi(d_i/2)l_i$, where $d_i$ is the diameter (width) and $l_i$ is the length for a skeleton pixel $i$. |
| 16/17/18 | VMeV/VmV/VMV | Median/min/max edge volume. The edge volume (in mm$^3$) corresponds to the volume of the same cylinder as in surface area, and is computed as $\sum_i V_i = \pi(d_i/2)^2 l_i$. |
| 19/20/21 | VMeO/VmO/VMO | Median/min/max edge orientation. The orientation is the angle (in the range $[-90°, 90°]$) between the x-axis and the major axis of the ellipse with the same second moments as the vein. |
| 22 | AN | Total number of areoles in the image patch. |
| 23/24/25 | AMeP/AmP/AMP | Median/min/max areole perimeter. The perimeter (in mm) is the distance along the pixels of the border of the areole. |
| 26/27/28 | AMeA/AmA/AMA | Median/min/max areole area. The areole area (in mm$^2$) is the number of pixels in each areole times the area of one pixel. |
| 29/30/31 | AMeCA/AmCA/ AMCA | Median/min/max areole convex area. The convex area (in mm$^2$) is the area of the convex hull for the areole. |
| 32/33/34 | AMeS/AmS/AMS | Median/min/max areole solidity. The solidity is a dimensionless parameter between 0 and 1 which measures the proportion of the pixels in the convex hull that are also in the area (ratio between the areole area and the convex area). |
| 35/36/37 | AMeMaA/ AmMaA/AMMaA | Median/min/max areole major axis. The major axis (in mm) corresponds to the ellipse with the same normalized second moments as the areole. |
| 38/39/40 | AMeMiA/AmMiA/ AMMiA | Median/min/max areole minor axis. The minor axis (in mm) corresponds to the ellipse with the same normalized second moments as the areole. |
| 41/42/43 | AMeE/AmE/AME | Median/min/max areole eccentricity. The eccentricity is a dimensionless parameter between 0 (a circle) and 1 (a line), which measures the ratio of the distance between the foci of the ellipse having the same normalized second moments as the areole and its major axis. |
| 44/45/46 | AMeEq/AmEq/ AMEq | Median/min/max areole equivalent diameter. The equivalent diameter (in mm) is the diameter of a circle having the same area as the areole. |
| 47/48/49 | AMeMD/AmMD/ AMMD | Median/min/max areole mean distance. The mean distance (in mm) is the mean value of the Euclidean distances between each areole pixel and the nearest vein pixel. |
| 50/51/52 | AMeVD/AmVD/ AMVD | Median/min/max areole variance distance. The variance distance (in mm) is the variance of the Euclidean distances between each areole pixel and the nearest vein pixel. |

state-of-the-art classifiers and are also very fast. Random Forests also allows to estimate the importance of input variables (in their original dimensional space).

The algorithm constructs a set of unpruned trees from $B$ random samples with replacement (bootstrap versions) of the original training dataset. For each random forest tree, a random sample of $m$ variables from the full set of $p$ variables ($m \leqslant p$) is selected to split the data at each node and grow the decision tree. The final classification result is the class corresponding to the majority vote of the ensemble of trees. In this work, we used 500 trees and a standard value of $m = \sqrt{p}$.

Random Forests has an internal procedure to estimate the relevance of the features. After training, the features are shuffled one at a time. An out-of-bag estimation of the prediction error is made on this permuted dataset. Intuitively, a feature which is not important to the model will not alter significantly the classification performance when shuffled. On the other hand, if the model made strong use of a certain feature, changing its values will produce an important decrease in performance. The relative loss in performance between the original dataset and the shuffled dataset is therefore related to the relative importance of the feature affected by the process.

### 2.5.2. Support Vector Machines (SVM)

Support Vector Machines (SVM) (Vapnik, 1995) is a state-of-the-art classifier which assumes that applying an appropriate

nonlinear mapping of the data into a sufficiently high dimensional space, two classes can be separated by an optimum hyperplane. This decision hyperplane is chosen in such a way that the distance between the nearest patterns of different classes (i.e., the margin) is maximized. SVM depends on a regularization parameter, $C$, which controls the trade-off between the complexity of the classifier and the number of allowed misclassifications. Inner validation was used in this work to set this parameter during the training phase.

The decision surface may be linear or nonlinear. In the latter case, a kernel function can be used to map the patterns into a high dimensional space. In this work, we considered SVMs with a Gaussian kernel (SVMG). The Gaussian standard deviation was optimized in a validation step during the training.

The feature ranking is performed by following the sensitivity analysis described in the paper by Guyon et al. (2002) for non-linear kernels, by sorting decreasingly the features according to the change they produce in the classification cost function when they are individually removed.

### 2.5.3. Penalized Discriminant Analysis (PDA)

Fisher's Linear Discriminant Analysis (LDA) (Hastie, Tibshirani, & Friedman, 2009) is a classical classifier and dimension reduction tool which searches for linear combinations of the features in such a way that the class means of the linear combinations are maximally separated relative to the intra-class variance. The classification of new observations is then performed by assigning them to the closest centroid according to a distance metric (typically the Mahalanobis distance) in the transformed space.

In order to improve LDA, Penalized Discriminant Analysis (PDA) was proposed by Hastie et al. (1995). PDA is a regularized version of LDA, which adds a penalty term to the intra-class covariance matrix. PDA is useful for image classification problems with large number of highly correlated features.

In this work, standard Ridge Regression (GenRidge) (Hastie et al., 2009) was used, which has the ridge constant $\lambda$ as the only free parameter. This constant penalizes high values of the fitted variables, and is similar to the $C$ parameter in SVM. This parameter was automatically selected using a validation set in the training phase.

The importance of each feature can be computed as the sum of the corresponding loadings (in absolute values) across the eigenvectors belonging to the largest eigenvalues (Song, Mei, & Li, 2010). In our case, we have only two eigenvectors, and we chose to keep both in order to analyze the feature relevance.

### 2.5.4. Recursive Feature Elimination (RFE)

Recursive Feature Elimination (RFE) was initially proposed in the work by Guyon et al. (2002), where it was implemented in conjunction with SVMs to perform feature selection. It is a kind of backward feature elimination algorithm (Kohavi & John, 1997) where the relevance of the features is evaluated by ranking subsets of features instead of ranking the features individually. In this way, features that are not relevant when considered alone may become important by complementing another features.

Even though Guyon et al. (2002) originally used SVMs in combination with RFE, any other classification algorithm can be used instead. In this work, we combined RFE with the three classifiers under consideration, i.e., RF, SVMG and PDA.

On a validation step, the procedure starts by training the classifier with the whole set of features, and ranking them according to the importance determined by the classifier. Iteratively, a subset with the lowest-ranked features found by the classifier is removed (we selected 10% of the current total features). The accuracy is computed at each iteration, corresponding to the accuracy achieved by using the remaining subset. Five-fold cross-validation

is used to determine the optimum cardinality of the best subset of features.

Finally, the whole training set is used to train the classifier using, initially, the entire set of features, and progressively eliminating the lowest-ranked subset of features until the previously determined optimum cardinality is reached. Then, this trained classifier is used to make predictions on an unseen test set.

## 3. Results and discussion

The total number of features computed per leaf rises to 208, i.e., 52 features × 4 patches (combined veins and 3 scales). As a preprocessing step, all the features exhibiting near zero variance across the examples were discarded. Also, the data were normalized (centered and scaled). For each one of the three classifiers described in Section 2.5, both the whole set of features and a subset composed by the optimal number of relevant features (according to RFE) were considered. We also compared the results obtained by using both the combined veins + 3 individual scales with the ones achieved by considering the features measured on the combined veins only. Classification was performed in R resorting to the `randomForest` (Liaw & Wiener, 2002), `mda`[1] and `e1071` (Meyer, 2009) packages.

First, we report the results for the recognition of the three different legume species, i.e., soybean, red bean and white bean. Next, we discuss the results for the soybean cultivar identification. In all the cases, we performed 10 runs of 10-fold cross validation to estimate the final accuracy of the classification procedure. We used 5-fold cross validation for the optimization of the parameters corresponding to each classifier as well as the best number of features for RFE.

### 3.1. Legume species recognition

The vein segmentation results are shown in Fig. 3„ top panel, for a soybean leaf, as well as the 100 × 100-pixel central patches used for feature extraction. Fig. 3(b) corresponds to the segmentation of the combined veins image. Fig. 3(c) is the 100 × 100-pixel central patch obtained from Fig. 3(b). Fig. 3(d)–(f) are the central patches cropped from the veins at scales #1 to #3, respectively. As it can be seen from this figure, scale #1 contains mainly the primary order veins while scale #2 preserves more detail about smaller veins. The veins at scale #1 are much thinner than at scale #2. Scale #3 is a less noisy version of scale #2, and also shows thicker veins. The combination of the veins at different scales (Fig. 3(b) and (c)) is the most complete result, providing detail on both primary and smaller veins. However, higher order veins (e.g., terminal veins) are not possible to segment since they are not visible (the images were scanned without any clearing or amplification procedures, as explained in Section 2.1). A similar analysis can be done for the segmentation of a white bean leaf (Fig. 3, middle panel) and a red bean leaf (Fig. 3, bottom panel).

The total accuracies and accuracies per class obtained by each classification algorithm are reported in Table 2 as mean ± standard error ($S_E$). For each alternative automatic classifier, namely RF, SVMG and PDA, the classification was performed by using 1) the whole set of features; and 2) only the subset with the most relevant features after performing RFE (as described in Section 2.5.4 and the introductory part of Section 3), namely RF RFE, SVMG RFE and PDA RFE. The first part of Table 2 shows the results of using the features measured on the combined veins only (52 features), whereas the second part of the table presents the results of using the combined veins and the 3 individual scales (208 features).

---

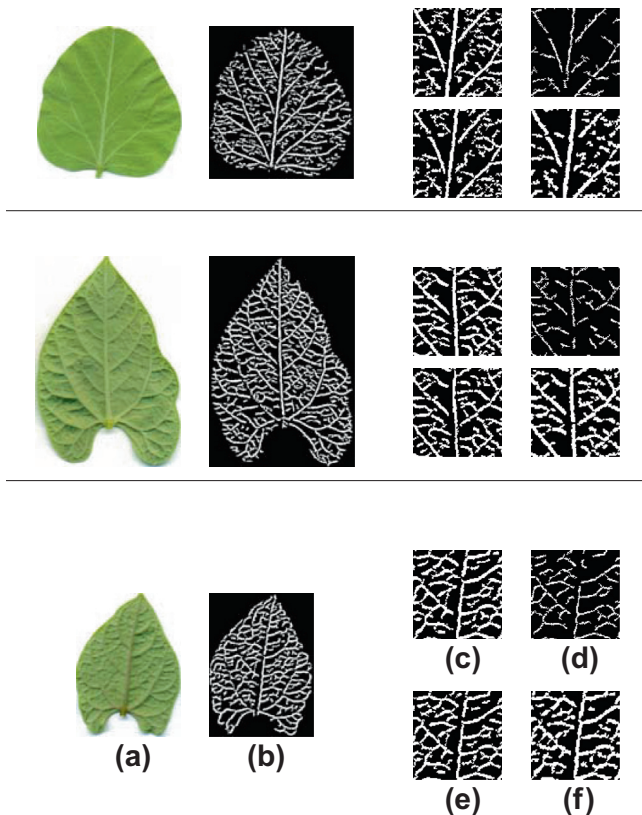[1] http://cran.r-project.org/web/packages/mda/index.html.

**Fig. 3.** Vein segmentation for a soybean leaf (top panel), a white bean leaf (middle panel) and a red bean leaf (bottom panel). (a) Original image, (b) combined veins, (c) central patch extracted from (b), (d) central patch extracted from scale #1, (e) central patch extracted from scale #2, and (f) central patch extracted from scale #3. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

As it can be observed from Table 2, the recognition of soybean leaves seems to be the easiest problem, since all the automatic classifiers configurations provide a mean accuracy over 95.5%. The recognition of the red bean leaves is also very good (mean accuracy over 84%). The classification of white bean leaves seems to be the most difficult, although the results are quite satisfactory (mean accuracy over 72%).

The accuracies achieved after performing feature selection via RFE are very similar to the ones obtained by using all the features,

both for the combined veins and the combined veins + 3 scales traits. The standard errors have very little differences, too. On the other side, the addition of the features measured on the 3 individual scales clearly increases the accuracies for all the classes and all the classifiers (with and without RFE) versus their counterparts using only the combined veins traits. These results highlight the benefit of considering the multiscale information extracted from the veins, showing the improvements introduced by the proposed approach.

We also analyzed the classification accuracies obtained by five human experts who developed manual classification of the same vein image patches. We found that the results of their classification were much more variable in comparison to the automatic classifiers. We report the average accuracies obtained by the experts considering only the patches from the combined veins as well as the average accuracies obtained by using (additionally) the patches at the 3 individual scales. It is evident that the 3 scales provide extra information to the experts which highly improves the performance of manual classification for white and red beans. For soybean, four experts out of the five achieve a performance which is similar or superior to the combined veins features case. Only one expert diminishes its performance when using the 3 scales, causing a small decrease in the mean value. However, it is a small loss against the important improvement achieved for the red and white bean classes. Overall, for the three species the total accuracy increases when using both combined veins and individual scales for manual classification.

From Table 2 it is clear that all of the automatic classifiers (either using feature selection or not) outperform manual classification for the three considered classes, apart from providing obvious advantages in repeatability, reliability and economy.

The results reported in Table 2 using the combined veins features only and without RFE (i.e., RF, SVMG and PDA) are consistent with previous results in the recent literature (Larese, Craviotto, Arango, Gallo, & Granitto, 2012; Larese et al., 2014). However, the usage of the combined veins + 3 scales features proposed in the present work outperform these previous results for all the cases.

We compared the usage of venation features versus the shape and texture features described in the work by Golzarian and Frick (2011). When following the approach of shape and texture features, we obtained an average total accuracy of ≈82% for legume classification with any of the 3 considered classifiers (RF, SVMG and PDA), which is much lower than the one we report by using venation features. Confusion matrices indicate that this result is

**Table 2**
Accuracy (mean $\pm S_E$) for legume species detection using 10 runs of 10-fold cross-validation.

| Classification algorithm | Per class accuracy (mean $\pm S_E$)% | | | Total accuracy (mean $\pm S_E$)% |
|---|---|---|---|---|
| | White bean | Red bean | Soybean | |
| *Combined veins* | | | | |
| RF | 72.26 ± 1.06 | 84.72 ± 0.65 | 96.00 ± 0.31 | 87.75 ± 0.32 |
| RF RFE | 72.30 ± 1.17 | 84.46 ± 0.66 | 95.62 ± 0.33 | 87.48 ± 0.36 |
| SVMG | 77.62 ± 1.00 | 85.68 ± 0.59 | 97.25 ± 0.25 | 89.72 ± 0.28 |
| SVMG RFE | 75.06 ± 0.97 | 86.22 ± 0.61 | 97.44 ± 0.24 | 89.48 ± 0.27 |
| PDA | 82.69 ± 0.88 | 85.83 ± 0.63 | 96.47 ± 0.27 | 90.39 ± 0.27 |
| PDA RFE | 81.87 ± 0.96 | 85.64 ± 0.63 | 96.35 ± 0.28 | 90.12 ± 0.30 |
| Manual classification | 66.43 ± 5.36 | 69.44 ± 6.79 | 98.29 ± 0.79 | 82.90 ± 1.62 |
| *Combined veins + 3 scales* | | | | |
| RF | 80.43 ± 0.97 | 89.03 ± 0.58 | 98.81 ± 0.17 | 92.09 ± 0.29 |
| RF RFE | 80.20 ± 0.92 | 89.29 ± 0.59 | 98.81 ± 0.16 | 92.13 ± 0.28 |
| SVMG | 81.33 ± 0.94 | 89.51 ± 0.65 | 98.15 ± 0.18 | 92.10 ± 0.28 |
| SVMG RFE | 81.02 ± 0.85 | 89.01 ± 0.67 | 98.15 ± 0.21 | 91.88 ± 0.28 |
| PDA | 90.92 ± 0.59 | 91.68 ± 0.51 | 98.98 ± 0.14 | 95.09 ± 0.21 |
| PDA RFE | 90.03 ± 0.61 | 90.32 ± 0.54 | 98.86 ± 0.14 | 94.43 ± 0.23 |
| Manual classification | 70.82 ± 13.15 | 83.28 ± 3.71 | 96.65 ± 0.85 | 87.32 ± 1.96 |

due to the mistakes in classification of white beans which are misclassified as red beans, since their leaves look very similar in shape and texture. However, multiscale venation differences are stronger and provide a better discrimination for this particular problem.

Fig. 4 depicts the distributions of the generalization errors obtained by each classifier over the 10 runs of 10-fold cross validation. The mean values of the manual classification accuracies using the features from the combined veins and combined veins + 3 scales are also included for comparison purposes.

Evidently, the distributions of RF, SVMG and PDA show an improvement in the performance when using the features from the combined veins + 3 individual scales. This can also be noticed for the manual classification. The distributions of PDA and PDA RFE with 208 features are completely over the manual classification accuracy with 208 features, and do not present outliers. All the distributions are approximately symmetric and have similar dispersion, except for RF and RF RFE with 52 features, for which it is slightly higher.

### 3.1.1. RFE analysis

Table 3 shows the 10 features most frequently selected as highly relevant for each algorithm along the 10 runs of 10-fold cross validation. It is noticeable that some of the features considered as relevant for each classifier are correlated, according to Larese et al. (2014). This is the case, for example, for AMA and AMCA selected by PDA RFE and SVMG RFE with 52 features. From Table 3 it can also be noticed that the only feature considered simultaneously relevant by all the algorithms is VMeW (highlighted in bold).

For the three classifiers with 208 features, the most frequently selected features are chosen from the combined veins and the 3 individual scales, showing the usefulness of these scales. The features which are considered highly relevant simultaneously by the three algorithms are VNN1, VNE1 and VMeW2 (highlighted in bold in Table 3).

### 3.2. Soybean cultivar recognition

Next we applied the proposed procedure to perform the recognition of three different soybean cultivars. Three exemplars (one per each class) are shown in the three panels of Fig. 5.

From these figures, an analysis similar to the one performed in Section 3.1 can be made regarding to the differences and

similarities between the three segmented individual scales and the combined veins patches. However, in this case there are not evident differences between the veins of the three cultivars. Moreover, when analyzing the dataset we noticed that there exists a high variability between individuals from the same class. Thus, this application problem is characterized by a relatively low inter-class and high intra-class variabilities.

We report in Table 4 the total and per class classification accuracies obtained by the different proposed automatic classifiers using both the whole set of features and feature selection by means of RFE. Since the current problem is more challenging than the species recognition one, it is expected to get a lower performance than for the species classification. However, we highly improved the total accuracies achieved by manual classification. The total accuracies are over 55% for all the algorithms under consideration. From Table 4 it is also noticeable that the recognition of cultivar 3 is the most difficult to achieve automatically (the best result is 34.28% for PDA RFE with 208 features). Similarly, human experts also have a very low performance, obtaining 39.47% and 43.98% of accuracy with the combined veins and the combined veins + 3 scales, respectively. However, all the automatic algorithms highly improve the manual classification both for cultivars 1 and 2. PDA RFE with 208 features presents a slightly higher average accuracy. All the algorithms have similar difficulties in separating the three cultivars, even though they all highly outperform the manual classification providing also with less variability.

The distribution of performances along the 10 runs of 10-fold cross validation for the different automatic classifiers (with and without feature selection) can be analyzed with the help of Fig. 6, where the total accuracy distributions are depicted for each algorithm. From this figure, it is evident that all the distributions are similar and almost symmetric, except for RF with 52 features (Fig. 6(a)) and PDA with 208 features (Fig. 6(b)), which present an asymmetry to the highest accuracies. The distribution of PDA RFE is entirely above manual classification both for 52 and 208 features, reaching higher accuracies in the last case. On the other hand, the distribution of accuracies for RF and RF RFE is partially below the manual classification with both numbers of features.

If the leaves to be classified are known to belong all to the same cultivar (as in the case of an unidentified seed lot), the accuracy of the system can be substantially improved by taking several exemplars from the same seed lot. The seed lot is then assigned to the most voted class.
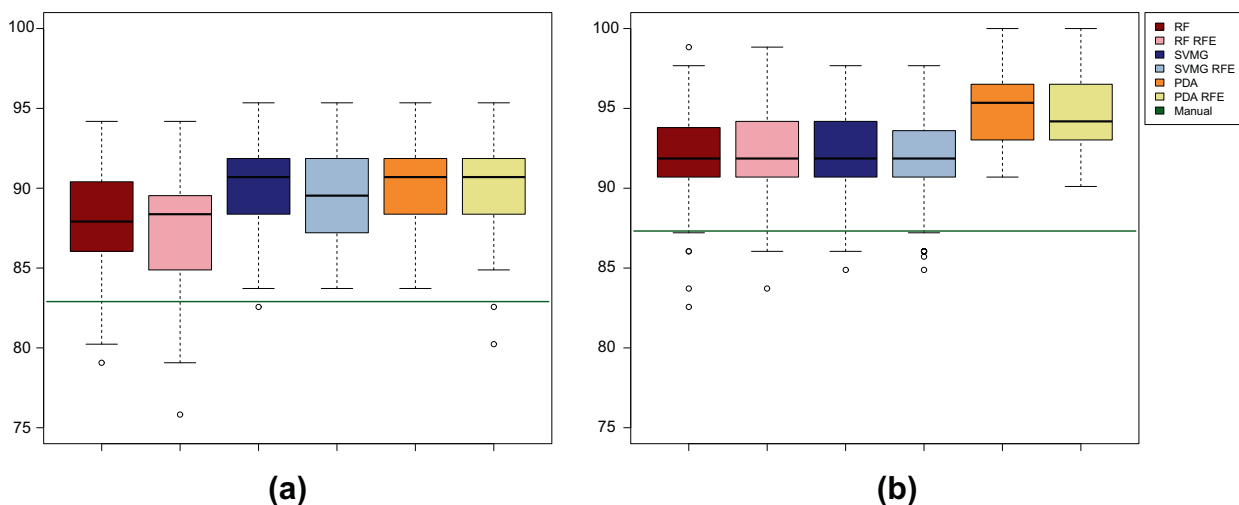


**Fig. 4.** Total accuracy distributions over 10 runs of 10-fold cross validation for the different classification algorithms (legume recognition problem). (a) 52 features and (b) 208 features.

**Table 3**
List of the 10 most selected features for each classifier and the legume classification problem. The percentage of times that each feature was selected is also shown.

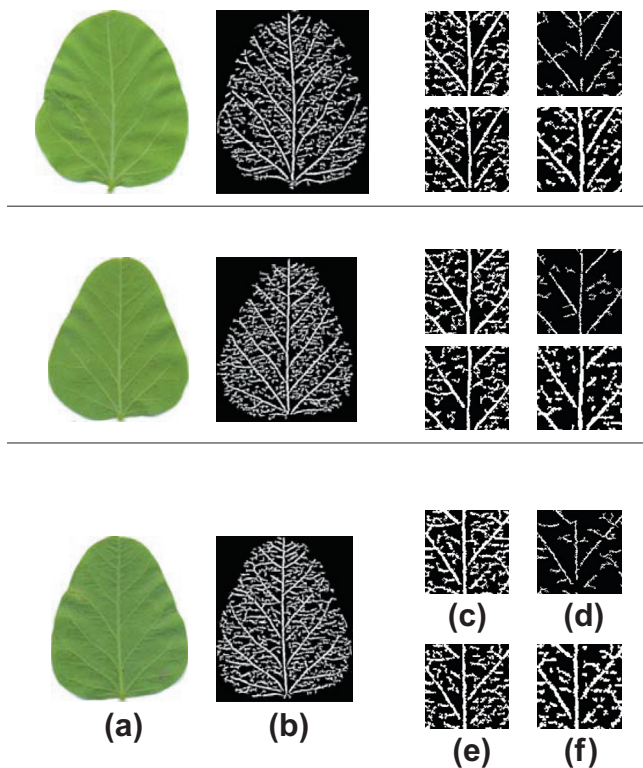| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| *RF with 52 features* | | | | | | | | | |
| VMeW | VMeV | VMeO | VMV | AN | AmS | VTNL | AMeMaA | AMeP | AMMiA |
| 100% | 100% | 100% | 100% | 100% | 100% | 99% | 99% | 93% | 91% |
| *RF with 208 features* | | | | | | | | | |
| VNE1 | VNN1 | VTNL1 | AMP1 | AMMD1 | AMVD1 | VMeW2 | AmS2 | AMEq2 | VMeW |
| 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% |
| *SVMG with 52 features* | | | | | | | | | |
| VNE | VMeW | VMeV | AMCA | VNN | AN | AmS | AMA | VTNL | VMW |
| 100% | 100% | 100% | 99% | 98% | 96% | 96% | 93% | 91% | 81% |
| *SVMG with 208 features* | | | | | | | | | |
| VNE1 | VMeW2 | VMeW | AMP1 | AMCA | VNN1 | AN | AMeMD | VNN3 | VMeL1 |
| 100% | 100% | 99% | 98% | 98% | 94% | 94% | 93% | 92% | 91% |
| *PDA with 52 features* | | | | | | | | | |
| VNE | VNN | VTNL | VMeW | VMV | AMA | AMCA | VML | AMeMD | AMEq |
| 100% | 100% | 100% | 100% | 100% | 100% | 100% | 97% | 97% | 97% |
| *PDA with 208 features* | | | | | | | | | |
| VNE1 | VNN1 | VTNL1 | AMCA1 | VMeW2 | VNE | VNN | VMeW | AMeMD | AMCA |
| 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% |



**Fig. 5.** Vein segmentation for a soybean leaf from cultivar 1 (top panel), cultivar 2 (middle panel) and cultivar 3 (bottom panel). (a) Original image, (b) combined veins, (c) central patch extracted from (b), (d) central patch extracted from scale #1, (e) central patch extracted from scale #2, and (f) central patch extracted from scale #3.

Assuming that all the classes are equally probable, and that the distribution of the classifier errors is uniform among the classes, the identification accuracy can be computed resorting to the binomial distribution $p' = (n!/((n/c)!(n-n/c)!)p^n(1-p)^{(n-n/c)}$, where $n$ is the number of exemplars to be classified, $c$ is the number of possible classes, $p$ is the average accuracy of the classifier (slightly higher than 50% is required at least) and $p'$ stands for the desired accuracy. For example, in order to achieve $p' = 99\%$ of classification accuracy for the PDA RFE with the combined veins + 3 scales features, from Table 4 we could take $p = (60 - 2S_E)\%$, resulting in

$n = 8$ leaves to be evaluated. The class is selected according to which the majority of the specimens were classified into. The lower the performance of the classifier, the higher number of exemplars ($n$) will be required to ensure a desired accuracy.

### 3.2.1. RFE analysis

We analyzed the number of times that each feature was selected as relevant by the classifiers and RFE along the 10 runs of 10-fold cross validation. In Table 5 we show the 10 features that were selected more frequently, as well as the percentage of times that they were selected. From this table, it can be noticed that VTNL is found relevant for all the considered classifiers (highlighted in bold). When using 208 features, the set of the 10 most frequently selected include features from the three individual scales and the combined veins image for all the classifiers. This reinforces the evidence that the 3 individual scales add important information to the problem, and that the three employed classifiers can take advantage of them in order to improve the classification. In addition, it can be noticed that, in this case, VTNL1 is considered relevant simultaneously by all the classifiers (highlighted in bold in Table 5).

We develop an automatic low cost procedure to classify legume varieties, based exclusively on the multiscale vein feature analysis of scanned leaf images. This method is useful when leaf shape, color and texture do not differ between the classes. We use state-of-the-art classifiers and feature selection techniques. Our method improves the previous results published in the recent literature concerning legume species recognition. Also, our method is tested on the yet untackled and more difficult problem of legume cultivars classification, where the leaves have all similar appearance. In both problems our method outperforms the human expert classification accuracy. However, in the second problem our accuracies are still low, even though they can be raised in the presence of seed lots where several leaf exemplars can be examined in the knowledge that they all belong to the same unknown class. In this case, the seed lot can be labeled according to the most frequent class.

Our method assumes that all the leaves were scanned in the same position, thus avoiding rotation dependence of the measured features. If this is not the case, a preprocessing step should be added to previously correct the position of all the leaves by checking the primary vein direction.

The method has been tested on a database composed by only three legume species, namely soybean, red bean and white bean, and only three soybean cultivars. It is necessary to test the method on an augmented database including new species and cultivars.

**Table 4**
Accuracy (mean $\pm S_E$) for soybean cultivars detection.

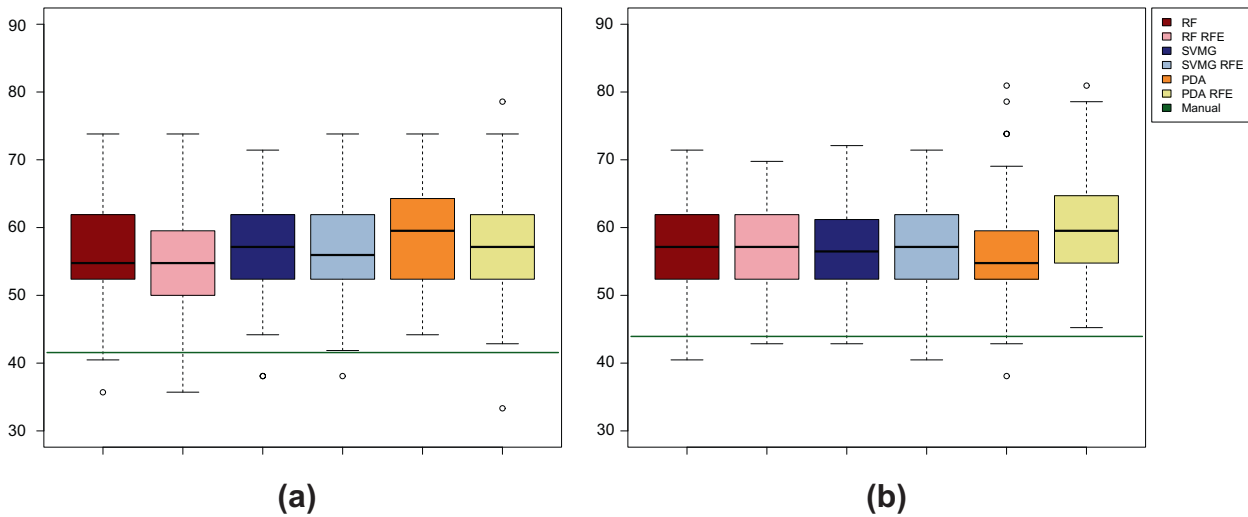| Classification algorithm | Per class accuracy (mean $\pm S_E$%) | | | Total accuracy (mean $\pm S_E$%) |
|---|---|---|---|---|
| | Cultivar #1 | Cultivar #2 | Cultivar #3 | |
| *Combined veins* | | | | |
| RF | 67.83 ± 1.19 | 54.18 ± 1.17 | 19.15 ± 1.64 | 56.66 ± 0.76 |
| RF RFE | 66.87 ± 1.20 | 52.02 ± 1.09 | 16.80 ± 1.60 | 55.04 ± 0.67 |
| SVMG | 70.92 ± 1.17 | 54.64 ± 1.15 | 5.78 ± 1.14 | 56.78 ± 0.69 |
| SVMG RFE | 70.26 ± 1.31 | 54.74 ± 1.19 | 5.62 ± 1.15 | 56.47 ± 0.73 |
| PDA | 69.49 ± 1.24 | 57.24 ± 1.23 | 19.32 ± 1.90 | 58.76 ± 0.74 |
| PDA RFE | 70.11 ± 1.23 | 53.24 ± 1.11 | 19.02 ± 1.78 | 57.31 ± 0.77 |
| Manual classification | 47.23 ± 8.06 | 35.76 ± 3.71 | 39.47 ± 4.32 | 41.56 ± 3.01 |
| *Combined veins + 3 scales* | | | | |
| RF | 67.63 ± 1.04 | 57.12 ± 1.06 | 13.27 ± 1.63 | 57.10 ± 0.67 |
| RF RFE | 67.05 ± 1.04 | 56.09 ± 1.11 | 16.10 ± 1.71 | 56.69 ± 0.63 |
| SVMG | 71.41 ± 1.15 | 53.36 ± 1.15 | 6.67 ± 1.29 | 56.57 ± 0.65 |
| SVMG RFE | 68.54 ± 1.15 | 56.06 ± 1.23 | 11.80 ± 1.56 | 56.93 ± 0.67 |
| PDA | 64.53 ± 1.08 | 54.60 ± 1.16 | 27.03 ± 2.00 | 56.17 ± 0.70 |
| PDA RFE | 68.22 ± 1.12 | 58.40 ± 1.16 | 34.28 ± 2.37 | 60.20 ± 0.75 |
| Manual classification | 44.95 ± 2.00 | 42.78 ± 5.37 | 43.98 ± 6.97 | 43.94 ± 2.48 |



**Fig. 6.** Total accuracy distributions over 10 runs of 10-fold cross validation for the different classification algorithms (soybean cultivar recognition). (a) 52 features, (b) 208 features.

**Table 5**
List of the 10 most selected features for each classifier and the cultivar classification problem. The percentage of times that each feature was selected is also shown.

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| *RF with 52 features* | | | | | | | | | |
| VNN | VTNL | VMeV | VNE | AMA | AMCA | AMEq | VMeA | VMeSA | AN |
| 100% | 100% | 99% | 98% | 98% | 98% | 98% | 97% | 95% | 89% |
| *RF with 208 features* | | | | | | | | | |
| VNE3 | VNE | VNN | VTNL | AMCA2 | AMEq2 | VTNL3 | AMEq1 | VTNL1 | AN |
| 100% | 100% | 100% | 100% | 99% | 99% | 99% | 98% | 97% | 97% |
| *SVMG with 52 features* | | | | | | | | | |
| AN | VTNL | AMCA | AMP | AMMaA | AMeE | VMeO | VNN | AMEq | VMeL |
| 98% | 94% | 87% | 78% | 77% | 76% | 73% | 71% | 71% | 70% |
| *SVMG with 208 features* | | | | | | | | | |
| VTNL | VTNL1 | AMMaA3 | AN | VMW3 | VMeV2 | VMeA2 | VMeO2 | AME3 | VNE3 |
| 99% | 94% | 82% | 78% | 77% | 74% | 71% | 68% | 65% | 63% |
| *PDA with 52 features* | | | | | | | | | |
| VTNL | VNN | VNE | VML | VMV | AMCA | AMMiA | VMeO | VMeA | AN |
| 100% | 99% | 93% | 91% | 91% | 82% | 63% | 60% | 57% | 51% |
| *PDA with 208 features* | | | | | | | | | |
| VNE1 | VNN1 | VTNL1 | VNN2 | VNE | VTNL | VMV | VML | VNE2 | AMCA3 |
| 99% | 99% | 99% | 92% | 90% | 88% | 88% | 85% | 79% | 79% |

## 4. Conclusions

In this work we show how an automatic image analysis and machine learning system can be implemented to classify leaves from different species and varieties. This system can provide a reliable, repeatable and economical means of recognizing plants, outperforming manual expert classification. By focusing on vein features only, we attempt to deal with the problem of visually similar leaves from different cultivars or varieties, which otherwise require to be processed by expensive methods, such as DNA analysis.

We implemented three automatic classifiers, namely Random Forests, PDA and SVM with Gaussian kernel. We show that our method works well independently of the employed automatic classifier. We also demonstrate how the usage of multiscale vein features improves the classification accuracy in contrast to single scale features for the problems under analysis. Additionally, we show how the three classifiers can be used in combination with RFE to estimate the relevance of the features in order to highlight possible vein feature patterns for each species and cultivar.

Our proposed system was tested on two different problems. The first one is the easiest one, dealing with different legume species, and has already been considered in the recent literature. The veins of these species present some differences that are perceivable by the human experts, judging by the high accuracies they achieve on this problem. However, our automatic procedure improves both the human accuracy and the previous results reported in the literature.

We also tested the proposed method on the more challenging problem of classifying cultivars from the same species, where the leaves differences within the same cultivar are approximately of the same order as the visual differences between leaves from different cultivars. In this case, the performance of the proposed procedure is much lower than for the first recognition problem, although the achieved average accuracy outperforms human experts results.

Even though we obtain low accuracies in this second problem, they can be improved in the case of seed lot classification by taking and classifying several leaf exemplars and labeling the seed lot according to the most frequent class.

Once the proposed system has been trained on the species and varieties of interest, it could be used to recognize new leaf specimens directly by scanning the new leaf. The system automatically processes the leaf image, segments the veins and computes the whole set of features. After this, the trained classifiers use these inputs to predict the class of the plant.

Our results are encouraging, but further work is needed aimed at extending this study to new species and varieties, augmenting the plant database. Additionally, future work includes the addition of new features in order to improve the accuracy for the cultivar recognition. In this direction, semantic relations between vein branches can be considered. Evaluating the implementation of the proposed whole system to run in a mobile device is also of great interest in the near future.

## Acknowledgments

## References

Agarwal, G., Ling, H., Jacobs, D., Shirdhonkar, S., Kress, W., Russell, R., et al. (2006). First steps toward an electronic field guide for plants. *Taxon, Journal of the International Association for Plant Taxonomy, 55*, 597–610.

Bama, B. S., Valli, S. M., Raju, S., & Kumar, V. A. (2011). Content based leaf image retrieval (CBLIR) using shape, color and texture features. *Indian Journal of Computer Science and Engineering, 2*(2), 202–211.

Breiman, L. (2001). Random forests. *Machine Learning, 45*, 5–32.

Camargo Neto, J., Meyer, G. E., Jones, D. D., & Samal, A. K. (2006). Plant species identification using Elliptic Fourier leaf shape analysis. *Computers and Electronics in Agriculture, 50*, 121–134.

Chaki, J., & Parekh, R. (2012). Designing an automated system for plant leaf recognition. *International Journal of Advances in Engineering & Technology, 2*(1), 149–158.

Clarke, J., Barman, S., Remagnino, P., Bailey, K., Kirkup, D., Mayo, S., Wilkin, P. (2006). Venation pattern analysis of leaf images. In *Advances in visual computing. Lecture notes in computer science (ISVC2006)* (Vol. 4292, pp. 427–436).

Du, J.-X., Wang, X.-F., & Zhang, G.-J. (2007). Leaf shape based plant species recognition [Special issue on intelligent computing theory and methodology]. *Applied Mathematics and Computation, 185*(2), 883–893.

Du, J.-X., Zhai, C.-M., & Wang, Q.-P. (2013). Recognition of plant leaf image based on fractal dimension features. *Neurocomputing, 116*, 150–156.

Golzarian, M. R., & Frick, R. A. (2011). Classification of images of wheat, ryegrass and brome grass species at early growth stages using principal component analysis. *Plant Methods, 7*(28).

Guyon, I., Weston, S., Barnhill, S., & Vapnik, V. (2002). Gene selection for cancer classification using support vector machines. *Machine Learning, 46*(1–3), 389–422.

Hastie, T., Buja, A., & Tibshirani, R. (1995). Penalized discriminant analysis. *Annals of Statistics, 23*(1), 73–102.

Hastie, T., Tibshirani, R., & Friedman, J. (2009). *The elements of statistical learning* (2nd ed.). Springer.

Im, C., Nishida, H., Kunii, T. L. (1998). Recognizing plant species by leaf shapes – A case study of the Acer family. In *International conference on pattern recognition* (Vol. 2, p. 1171).

Kohavi, R., & John, G. H. (1997). Wrappers for feature subset selection. *Artificial Intelligence, 97*, 273–324.

Larese, M. G., Craviotto, R. M., Arango, M. R., Gallo, C., & Granitto, P. M. (2012). Legume identification by leaf vein images classification. In L. Alvarez, M. Mejail, L. Gomez, & J. Jacobo (Eds.), *Progress in pattern recognition, image analysis, computer vision, and applications. Lecture notes in computer science* (Vol. 7441, pp. 447–454). Berlin Heidelberg: Springer.

Larese, M. G., Namías, R., Craviotto, R. M., Arango, M. R., Gallo, C., & Granitto, P. M. (2014). Automatic classification of legumes using leaf vein image features. *Pattern Recognition, 47*(1), 158–168.

Liaw, A., & Wiener, M. (2002). Classification and regression by randomforest. *Rnews, 2*(3), 18–22.

Meyer, D. (2009). Support vector machines. The interface to libsvm in package e1071.

Park, J., Hwang, E., & Nam, Y. (2008). Utilizing venation features for efficient leaf image retrieval. *Journal of Systems and Software, 81*(1), 71–82.

Price, C. A., Symonova, O., Mileyko, Y., Hilley, T., & Weitz, J. S. (2011). Leaf extraction and analysis framework graphical user interface: Segmenting and analyzing the structure of leaf veins and areolas. *Plant Physiology, 155*, 236–245.

Pydipati, R., Burks, T. F., & Lee, W. S. (2006). Identification of citrus disease using color texture features and discriminant analysis. *Computers and Electronics in Agriculture, 52*, 49–59.

Sack, L., Dietrich, E. M., Streeter, C. M., Sanchez-Gomez, D., & Holbrook, N. M. (2008). Leaf palmate venation and vascular redundancy confer tolerance of hydraulic disruption. *Proceedings of the National Academy of Sciences of the United States of America, 105*, 1567–1572.

Scoffoni, C., Rawls, M., McKown, A. D., Cochard, H., & Sack, L. (2011). Decline of leaf hydraulic conductance with dehydration: Relationship to leaf size and venation architecture. *Plant Physiology, 156*, 832–843.

Soille, P. (1999). *Morphological image analysis: Principles and applications*. Springer-Verlag.

Solé-Casals, J., Travieso, C. M., Alonso, J. B., Ferrer, M. A. (2008). Improving a leaves automatic recognition process using PCA. In *IWPACBB* (pp. 243–251).

Song, F., Mei, D., Li, H. (2010). Feature selection based on linear discriminant analysis. In *International conference on intelligent system design and engineering application (ISDEA), 2010 Vol. 1* (pp. 746–749).

Sonka, M., Hlavac, V., & Boyle, R. (2008). *Image processing analysis and machine vision*. Thomson.

Umbaugh, S. E. (2005). *Computer imaging: Digital image analysis and processing*. CRC Press.

Valliammal, N., Geethalakshmi, S. (2011). Hybrid image segmentation algorithm for leaf recognition and characterization. In *International conference on process automation, control and computing (PACC), 2011* (pp. 1–6).

Vapnik, V. (1995). *The nature of statistical learning theory*. Springer-Verlag.