

Fisher Vectors for Leaf Image Classification: an Experimental Evaluation

Javier A. Redolfi¹, Jorge A. Sánchez² y Julián A. Pucheta⁴

¹ Centro de Investigación en Informática para la Ingeniería,
FRC, UTN, Córdoba, Argentina,
jredolfi@frc.utn.edu.ar

² CONICET, FaMAF, UNC, Córdoba, Argentina

³ Laboratorio de Matemática Aplicada al Control,
UNC, Córdoba, Argentina

Abstract. In this work we present an experimental analysis of the use of exponential family Fisher vector to solve the problem of visual plant identification. We make a comparison of the encoding of different descriptors with this framework and we evaluate the performance on public datasets and compare these results with state of the art methods proposed in the literature. We show that eFV framework performs very well in the problem of plant classification.

Keywords: plant identificacion, exponential family Fisher vectors, image classification

1 Introduction

In recent years, there has been an increasing interest on the problem of plant species classification using visual information [3, 7, 13, 20]. Some reasons of this are the large number of endangered species and the high rates of deforestation due to the shift of the agricultural frontier and a poor urban planning. Plants have a crucial role in the life on earth and their carelessness cause irreversible problems to our society, such as global warming, loss of biodiversity and environmental damage [3, 22]. The problem presents a very interesting challenge, because it is almost impossible for common people and very difficult for trained ones such as farmers, wood exploiters or even botanists [5]. The reasons of this are many, among which we can name the large number of species, accounted for approximately 200000, the vast intra-class variability and a high visual similarities between classes [20].

In this paper we address the problem using a recently proposed encoding called exponential family Fisher vectors (eFV) [19]. This encoding is a generalization of Gaussian based Fisher vector (FV) to the exponential family distribution, which allows to encode local descriptors in a large number of input domains such as real, integer or binary vectors and symmetric positive definite matrices (SPDM).

2 Related Work

There exists a large body of work on the leaf image classification problem. Most of which has been focused on the design of different preprocessing [24, 25] and feature extraction techniques [13, 22] as well as classification algorithms [9, 20] specifically designed for this problem. Regarding the feature extraction stage, methods can be grouped into two main categories: those using global features and those using local image descriptors.

In [24, 25], the authors propose to use shape and texture global features obtained after a segmentation step for the classification of leaves images. In [13] a system using geometric descriptors, multi-scale distance matrix, invariant moments and a new set of global descriptors is proposed. The computation of these descriptors requires a contour extraction step which, according to the authors, accounts for one of its main limitation. In [20] the authors propose a semi-automatic approach based on global features that requires the user to mark the base and apex of the leaf. The method presented in [3] is based on a set of global descriptors which are sensitive to rotations of the image. Therefore, an image alignment preprocessing step is mandatory prior to feature extraction.

With respect to the local descriptors based methods, in [9] is proposed a system based on sparse coding of SIFT descriptors and a similar method using a combination of descriptors including SIFT is presented in [16]. In [1], the authors propose the use of different local descriptors (SURF, Fourier, Rotation Invariant, LBP) encoded with FV to classify images of leaves taken with a natural background. In that work descriptors are calculated over Harris interest points and classified with an SVM in a OvA configuration. The authors of [15] use local descriptors (4 versions of SIFT and self-similarity) augmented with a polynomial method that takes into account neighbors descriptors and then encoded with FV. In [4] FV over SIFT and color moments are combined with CNN (Convolutional Neural Networks), including a preprocessing step to get the most representative bounding box of the image.

In this work we propose the use of local descriptors, encoded with a recent proposed framework, termed exponential family Fisher vectors (eFV) [19].

3 Method Description

The proposed pipeline contains three stages, the first is dense extraction of visual descriptors, then these descriptors may or may not be reduced in dimensionality with PCA, after that, encoding is performed with eFV and finally these vectors are classified using SVM. A diagram of the pipeline is shown in figure 1 and in the following we explain each of the parts.

3.1 Descriptors

Descriptors are extracted densely on a regular grid with the same step in both directions. Furthermore, these are calculated in the original image and in four

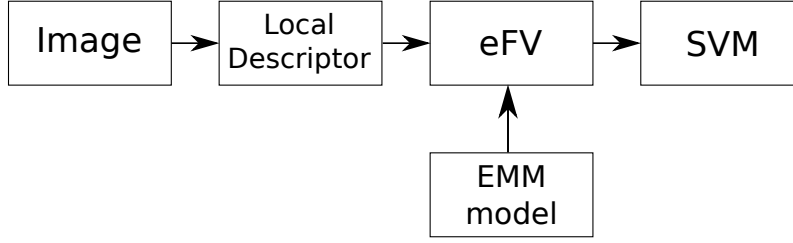


Fig. 1. Block diagram of the proposed system.

scales, with a scaling factor of $\frac{1}{\sqrt{2}}$. This value is the normal choice for this type of encoding [19].

The selected descriptors for this work are SIFT, binarized SIFT (BinSIFT), BRIEF, LBP and a variation of COV [21]. PCA was only applied to SIFT and BinSIFT descriptors.

3.2 Exponential Family Fisher Vector

The FV [18] representation is one of the most robust for image classification [14] and fine-grained classification [8]. This representation encodes an image as a gradient vector that characterizes the distribution of a set of low-level descriptors with respect to the parameters of a probabilistic generative model which in case of the traditional FV, corresponds to a mixture of multivariate Gaussian pdfs with diagonal covariances. The eFV generalizes the FV by considering mixtures of a more general class of distributions known as the exponential family. This allow the model to deal with input spaces other than \mathbb{R}^D in a principled manner. Next, we provide a brief overview of the eFV representation. More details can be found in [19].

Let $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, $\mathbf{x}_i \in \mathbb{R}^D$, a set of low-level descriptors extracted from image I . We model its generation process as a mixture distribution of the form:

$$P(\mathbf{X}|\lambda) = \prod_{i=1}^N \sum_{k=1}^K w_k p(\mathbf{x}_i|\eta_k), \quad w_k > 0 \quad \forall k, \quad \sum_{k=1}^K w_k = 1 \quad (1)$$

with $\lambda = \{w_k, \eta_k : k = 1, \dots, K\}$ and

$$p(\mathbf{x}|\eta_k) = h(\mathbf{x}) \exp(\langle \eta_k, T_k(\mathbf{x}) \rangle - \psi(\eta_k)) \quad (2)$$

is a member of the exponential family. Here, $T_k(\mathbf{x})$ is the vector of sufficient statistics, $\psi(\eta_k)$ is known as *partition function* and $h(\mathbf{x})$ is a normalizer which is independent of η_k . Given $P(\mathbf{X}|\lambda)$, the eFV of I is computed as the normalized gradient w.r.t λ of the log-likelihood of \mathbf{X} as:

$$g(\mathbf{X}) \triangleq L_\lambda \nabla_\lambda \log P(\mathbf{X}|\lambda) = \sum_{i=1}^N L_\lambda \nabla_\lambda \log \left(\sum_{k=1}^K w_k p(\mathbf{x}_i|\eta_k) \right) \quad (3)$$

Table 1. Descriptors and corresponding distribution used for encoding.

Descriptor	Input Domain	Distribution	$T_k(\mathbf{x})$	$\psi(\eta)$	$h(\mathbf{x})$
SIFT	\mathbb{R}^D	Gaussian	$(x_1, x_1^2, \dots, x_D, x_D^2)$	$-\frac{1}{3}\eta_1^T \eta_2^{-1} \eta_1 - \frac{1}{2} \log -2\eta_2 $	1
BRIEF, BinSIFT, LBP	$\{0, 1\}^D$	Bernoulli	\mathbf{x}	$\sum_{j=1}^D \log(1 + e^{\eta_j})$	1
COV	$D \times D$ SPDM	Wishart	\mathbf{x}	$\log \Gamma_n(\frac{n}{2}) - \frac{n}{2} \log \eta $	$ \eta ^{(n-D-1)/2}$

L_λ is a normalizer obtained from the Cholesky decomposition of the inverse of the Fisher information matrix of $P(\cdot|\lambda)$.

Table 1 show the descriptor and its corresponding exponential family distribution used for encoding.

3.3 Classifier

For eFV classification we used SVM with a linear kernel trained with SGD, because it is the normal selection for this type of codifications [18, 19]. The use of non-linear kernels is problematic due to the high dimensionality of the vectors.

4 Experiments

To evaluate the eFV encoding, we perform experiments on different public datasets, commonly used for this task and we compare our results with different state of the art algorithms.

4.1 Datasets

The first dataset is the presented in [23], known as Flavia, which contains 1907 images of leaves from 32 classes of trees, with a minimum of 50 samples per class and a maximum of 72. The normal procedure of evaluation is to leave 10 samples of each class for test and train on the rest.

The second dataset is known as Foliage [11], which contains 120 samples for each of 60 species of trees. The recommended procedure of evaluation is to take 100 samples for training and 20 for testing for each class.

The last two, are the used on the plant identification challenge organized in the ImageCLEF 2012 and 2013. The first of these datasets, PlantCLEF2012 [6], contains 11572 images of 126 species of trees divided in three types, scan, scan-like and photograph. The second, PlantCLEF2013 [5], contains 26077 images of 250 tree species of two types, sheet as background and natural background. Also, the NaturalBackground images are divided into 5 types, images of entire plant, flower, fruit, leaf and stem.

4.2 Experimental Configuration

As already mentioned, local descriptors were calculated on a regular grid and in four image scales with a factor of $\frac{1}{\sqrt{2}}$. In the case of SIFT and BinSIFT

Table 2. eFV configurations and short names.

Short Name	Descriptor	PCA	Exponential Mixture Model
BRIEF-BMM-eFV	BRIEF	No	Bernoulli
SIFT-PCA-GMM-eFV ⁴	SIFT	Yes	Gaussian
COV-WMM-eFV	Covariance	No	Wishart
LBP-BMM-eFV	LBP	No	Bernoulli
BinSIFT-BMM-eFV	BinSIFT	Yes	Bernoulli

descriptors, a step of dimensionality reduction using PCA was applied and the resulting dimensionality was 78. On these descriptors, a 256 component family exponential mixture model was fitted, which was then used to calculate the eFV encoding, according to the configuration shown in the table 1. Table 2 shows a resume of the different eFV configurations and its short name for further reference. eFV computing was done with the library mentioned in [19].

Furthermore, we propose the use of the results obtained using the descriptors based on CNN proposed in [17] as a baseline for comparison. In that work, the authors show that features obtained from CNN nets should be used as the first candidate in most visual recognition tasks. These descriptors were computed such as the output of the 7th layer (fc7) of the convolutional neural network available in [10] and then classified with an SVM. This baseline is referred in the following as CNN+SVM.

4.3 Results

Table 3 shows the accuracy of different configurations of the proposed method on the Flavia and Foliage datasets together with recent results available in the literature. The accuracy is obtained as the percent of well classified samples.

As can be seen from table 3 the best accuracy of eFV are obtained with SIFT and COV descriptors, and their accuracy on Flavia and Foliage is above the obtained with recent proposed methods in the literature. In the Foliage dataset the baseline CNN+SVM has the best accuracy.

In tables 4 and 5, we compare the results of our algorithm with the best results on the PlantCLEF2012 and PlantCLEF2013 challenges. The score is computed using the scripts provided with the datasets. In bold letters we highlight the best accuracy for each type of image. For these two datasets we only show the accuracy for SIFT and COV descriptors.

In PlantCLEF2012 dataset (table 4) the encoding of SIFT descriptors with eFV shows the best performance for Scan-like, Photos and Average, and the baseline system CNN+SVM, shows the best performance for Scan type of images.

For the dataset PlantCLEF2013, the best performance for SheetAsBackground images is obtained with the method proposed in [24] but this method fails for the NaturalBackground images as can be seen in table 5. The cause of

⁴ This configuration is the traditional FV [18].

Table 3. Accuracy of different configurations of eFV and results in the literature on datasets Flavia and Foliage

Method	Acc. Flavia	Acc. Foliage
CNN+SVM	99.06	99.33
SIFT-PCA-GMM-eFV	99.06	98.75
COV-WMM-eFV	99.38	98.25
LBP-BMM-eFV	95.62	93.25
BinSIFT-BMM-eFV	89.06	94.33
BRIEF-BMM-eFV	74.06	67.83
GLC [13]	93.00	-
SC [9]	95.47	-
CS [20]	97.00	-
GLS [12]	97.19	95.00
ICM [22]	97.82	-

Table 4. Classification results on PlantCLEF2012 for the 3 types of images and on average.

Method	Scan	Scan-like	Photos	Average
CNN+SVM	0.65	0.51	0.40	0.520
SIFT-PCA-GMM-eFV	0.62	0.74	0.44	0.60
COV-WMM-eFV	0.481	0.432	0.240	0.384
SABANCI OKAN [25]	0.58	0.55	0.22	0.16
INRIA [2]	0.39	0.59	0.21	0.40
LSIS DYNI [16]	0.41	0.42	0.32	0.42

this behavior is a preprocessing segmentation step which is inapplicable for NaturalBackground images. For this type of images, one of the best performance is achieved with the method presented in [15] based in a complex scheme of late-fusion of 4 versions of SIFT and self-similarity encoded with a polynomial embedding of descriptors encoded with FV. Also, the last method uses metadata information of the test set, in particular the type of NaturalBackground image, in contrast to our. Again, the baseline CNN+SVM has the best performance for one type of images.

Table 5. Classification results on PlantCLEF2013 for the 2 types of images.

Method	SheetAsBackground	NaturalBackground
CNN+SVM	0.557	0.403
SIFT-PCA-GMM-eFV	0.594	0.365
COV-WMM-eFV	0.363	0.181
SABANCI OKAN [24]	0.607	0.181
NlabUTokio [15]	0.502	0.393

5 Conclusions

We present a detailed empirical evaluation of different eFV configurations applied to the problem of plant identification. We perform experiments in different public datasets and compare our results with state of the art algorithms. The results in some experiments are better than the state of the art and in most of the cases the best eFV configuration is SIFT descriptors encoded with GMM based eFV. But the baseline using CNN and SVM also performs well and this performance can be explained by the power of discrimination of these descriptors.

The advantages of the proposed method are that it does not need a preprocessing step for the leaf contour extraction because it is based on local descriptors, it permits the use of different descriptors in an unified framework, it is not based in handcrafted or ad-hoc descriptors and it is simpler than some of the existing algorithms. Furthermore, unlike other methods it can be applied on images of leaves with a simple background or with complex background as we demonstrated on the experiments.

This research has shown that the encoding of SIFT descriptors with eFV is a good choice to solve plant identification problems. It was also shown that the complex CNN descriptors works well too.

References

1. Bakic, V., Mouine, S., Ouertani-Litayem, S., Verroust-Blondet, A., Yahiaoui, I., Goëau, H., Joly, A.: Inria's participation at ImageCLEF 2013 plant identification task. In: CLEF (Online Working Notes/Labs/Workshop) 2013 (2013)
2. Bakic, V., Yahiaoui, I., Mouine, S., Ouertani, S.L., Ouertani, W., Verroust-Blondet, A., Goëau, H., Joly, A.: Inria IMEDIA2's participation at ImageCLEF 2012 plant identification task. In: CLEF (Online Working Notes/Labs/Workshop) 2012 (2012)
3. Chaki, J., Parekh, R., Bhattacharya, S.: Plant leaf recognition using texture and shape features with neural classifiers. *Pattern Recognition Letters* 58, 61–68 (2015)
4. Chen, Q., Abedini, M., Garnavi, R., Liang, X.: IBM research Australia at LifeCLEF 2014: Plant identification task. In: Working notes of CLEF 2014 conference (2014)
5. Goëau, H., Bonnet, P., Joly, A., Bakic, V., Barthélémy, D., Boujema, N., Molino, J.F.: The ImageCLEF 2013 plant identification task. In: CLEF (2013)
6. Goëau, H., Bonnet, P., Joly, A., Yahiaoui, I., Bakic, V., Barthélémy, D., Boujema, N., Molino, J.F.: The ImageCLEF 2012 plant identification task. In: CLEF (2012)
7. Goëau, H., Joly, A., Bonnet, P., Selmi, S., Molino, J.F., Barthélémy, D., Boujema, N.: LifeCLEF plant identification task 2014. CLEF2014 Working Notes. Working Notes for CLEF 2014 Conference, Sheffield, UK, September 15-18, 2014 pp. 598–615 (2014)
8. Gosselin, P.H., Murray, N., Jégou, H., Perronnin, F.: Revisiting the Fisher vector for fine-grained classification. *Pattern Recognition Letters* 49, 92–98 (2014)
9. Hsiao, J.K., Kang, L.W., Chang, C.L., Lin, C.Y.: Comparative study of leaf image recognition with a novel learning-based approach. In: Science and Information Conference (SAI), 2014. pp. 389–393. IEEE (2014)

10. Jia, Y., Shelhamer, E., Donahue, J., Karayev, S., Long, J., Girshick, R., Guadarrama, S., Darrell, T.: Caffe: Convolutional architecture for fast feature embedding. In: Proceedings of the ACM International Conference on Multimedia. pp. 675–678. ACM (2014)
11. Kadir, A., Nugroho, L., Susanto, A., Santosa, P.: Neural Network Application on Foliage Plant Identification. *International Journal of Computer Applications* 29(9), 15–22 (2011)
12. Kadir, A.: A Model of Plant Identification System Using GLCM, Lacunarity And Shen Features. arXiv preprint arXiv:1410.0969 (2014)
13. Kalyoncu, C., Önsen Toygar: Geometric leaf classification. *Computer Vision and Image Understanding* 133, 102–109 (2015)
14. Ken Chatfield, Victor Lempitsky, A.V., Zisserman, A.: The devil is in the details: an evaluation of recent feature encoding methods. In: Proc. BMVC. pp. 76.1–76.12 (2011), <http://dx.doi.org/10.5244/C.25.76>
15. Nakayama, H.: Nlab-utokyo at ImageCLEF 2013 plant identification task. In: Working notes of CLEF 2013 conference (2013)
16. Paris, S., Halkias, X., Glotin, H.: Participation of LSIS/DYNI to ImageCLEF 2012 plant images classification task. In: CLEF (Online Working Notes/Labs/Workshop) (2012)
17. Razavian, A.S., Azizpour, H., Sullivan, J., Carlsson, S.: CNN Features off-the-shelf: an Astounding Baseline for Recognition. In: Computer Vision and Pattern Recognition Workshops (CVPRW), 2014 IEEE Conference on. pp. 512–519. IEEE (2014)
18. Sánchez, J., Perronnin, F., Mensink, T., Verbeek, J.J.: Image Classification with the Fisher Vector: Theory and Practice. *International Journal of Computer Vision* 105(3), 222–245 (2013)
19. Sánchez, J., Redolfi, J.: Exponential family Fisher vector for image classification. *Pattern Recognition Letters* 59, 26–32 (2015)
20. Sfar, A.R., Boujemaa, N., Geman, D.: Confidence Sets for Fine-Grained Categorization and Plant Species Identification. *International Journal of Computer Vision* pp. 1–21 (2014)
21. Tuzel, O., Porikli, F., Meer, P.: Region covariance: A fast descriptor for detection and classification. In: Computer Vision–ECCV 2006, pp. 589–600. Springer (2006)
22. Wang, Z., Sun, X., Ma, Y., Zhang, H., Ma, Y., Xie, W., Zhang, Y.: Plant recognition based on intersecting cortical model. In: Neural Networks (IJCNN), 2014 International Joint Conference on. pp. 975–980. IEEE (2014)
23. Wu, S.G., Bao, F.S., Xu, E.Y., Wang, Y.X., Chang, Y.F., Xiang, Q.L.: A leaf recognition algorithm for plant classification using probabilistic neural network. In: Signal Processing and Information Technology, 2007 IEEE International Symposium on. pp. 11–16. IEEE (2007)
24. Yanikoglu, B., Aptoula, E., Yildiran, S.T.: Sabanci-Okan system at ImageCLEF 2013 plant identification competition. In: Working notes of CLEF 2013 conference (2013)
25. Yanikoglu, B.A., Aptoula, E., Tirkaz, C.: Sabanci-Okan system at ImageClef 2012: Combining features and classifiers for plant identification. In: CLEF (Online Working Notes/Labs/Workshop) (2012)