

Plant Leaf Recognition

Albert Liu, albertpl@stanford.edu Yangming Huang, yangming@standford

I. INTRODUCTION

Automatic plant species recognition with image processing has application in weeds identification, species discovery, plant taxonomy, natural reserve park management and so on [3]. It is considered as a fine-grained image recognition problem which is hard to solve since

- 1) The subtle differences between different species in the same class. Sometime such fine differences can be even challenging to human experts.
- 2) Typically this requires large training data but it is not feasible due to the number of species (over 220000 [2]).



In this report, we describe our exploration with this problem, using traditional handcrafted features and features extracted from pretrained deep convolution neural network (ConvNets). The rest of the report is organized as follows. In section III, we describe the data set. Section.

II. RELATED WORK

Research on automatic leaf classification from image has been active since 2000. Lots of hand-crafted features have since proposed, ranging from shape based, to statistical texture and margin [2] [3] [1]. Also generic computer vision object recognition/detection features, such SIFT and HOG are studied for this problem. [TODO] Most of such manually engineered features achieve excellent accuracy on clean images taken in controlled conditions, which consist of one single well aligned leave on contrasting background, such as those images in data set [15]. Recently, with the huge success of ConvNets, particularly the winners of ILSVRC [11], researchers start to apply ConvNets to this problem [TODO]. [TODO] have suggested that generic features can be extracted from large ConvNet and yield very good results on fine-grained classification problems even without fine-tuning the pretrained model.

III. DATASET



We found two types of data set

- 1) Clean images, which consists of well aligned leaf on single contrasting background, with little or no variations of luminance or color.

Name	Species	Samples Per Species
Swedish [15]	15	75
Flavia [16]	33	~ 60

- 2) ImageCLEF [17], which is collected through crowd sourced application. This is a much more noisy datasets with variations on lighting conditions, viewpoints, background clutters and even occlusions. The dataset can be further split into two subset: uniform, which is taken in a more controlled environment, and natural-background, which is taken in a natural environment.

Name	Species	train samples	test samples
uniform	66	9607	1194
natural	57	2585	521

IV. APPROACH

A. Overview

Here is the pipeline of our system.

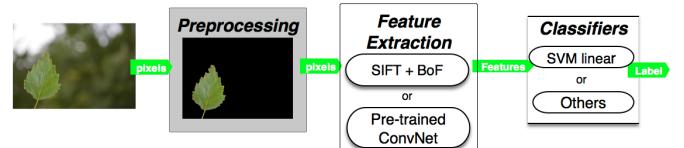


Fig. 1. Overview of the system

Firstly, during preprocessing, we apply CLAHE[27] to reduce lighting condition variation and then resize raw images to fit the next layer. We also attempt to remove background with the following techniques in a heuristic way

- 1) use K-means ($K=2$) and throw away the background
- 2) find the convex hull containing the largest N contours and then use GrabCut[26] to segment leaf out of the background clutter.

Next we extract features

- 1) ConvNets. We take transfer learning approach, in order to make use of the power of ConvNets under the constrains of time and computations. Specifically, we take a couple of ConvNets that are pretrained on ImageNet for ILSVRC object classification task, remove top FC layers and then treat the rest of the ConvNet as fixed feature extractor. The CNN codes are our features. To battle overfitting issue, we augment our input images.
- 2) Traditional SIFT + Bag of Features Key points are densely sampled and SIFT feature descriptors are retrieved at each key point. We fix size of the codebook (K) as 1000/300 for different data sets.

Finally, we train a simple classifier from the feature vectors and then predict labels for our test data.

B. ConvNets approach

1) *Setup:* For ConvNets, we used Keras framework [18] with Tensorflow backend of GPU support: NVIDIA GeForce GT 750M 2048 MB. As an alternative, we also run on CPU given the limited graphic memory of our GPU which is crucial for a deeper ConvNets architecture.

2) *Exploration:* We started by training ConvNets classifier from scratch, following the guidelines below:

- 1) Convolutional layer learning features from general to specific, giving more layers helps with the transition. Study on deep convolutional nets suggests that deeper models are preferred under a parameter budget[24].
- 2) Dropout reduces overfitting [25] [11]
- 3) Use aggressive pooling to reduce the dimensionality.

We designed our ConvNets with the architecture below:

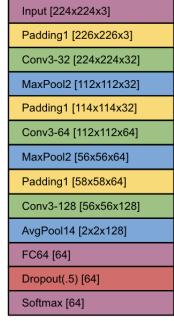


Fig. 2. Customized ConvNets Classifier Architecture

Using cross entropy to measure loss, we got the learning curve:

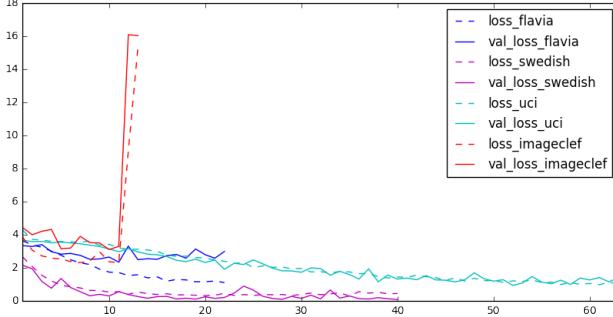


Fig. 3. Loss curve of ConvNets Classifier

The loss of swedish dataset [15] keeps decreasing as expected. But for flavia dataset [16], the training loss is keep decreasing, but at some point, validation loss stop changing.

For imageclef uniform dataset [17], it actually stopped learning and the loss goes back up and stay there.

The accuracy is consistent from the loss, where accuracy of swedish converges at 92%, with a test accuracy of 90.46%, but it doesn't work well with other datasets especially imageclef dataset.

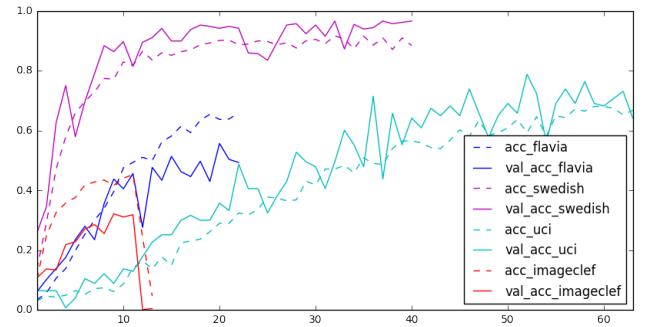


Fig. 4. Accuracy curve of ConvNets Classifier

Compare the datasets, intuitively, there are a few things to consider:

- The lighting variation and background clutter. Swedish has very clean background, where for ImageCLEF uniform image quality varies in a large range.
- More species with less samples per classes.

Some of our data set such as ImageCLEF2013 is even more challenging in term of clutter and occlusion, which we didn't cover in the customized ConvNets classifier experiment. A deeper ConvNets with more aggressive filters is needed to extract the features more efficiently and to deal with the noise and variation of datasets. But given the constrains of time and computations, we can't afford to train a deeper network with much more parameters.

As a common practice, we switched to use pre-trained weights of proven ConvNets architecture.

- 3) *Data Augmentation:* Before we start with the pre-trained weights, application of data augmentation is helpful.

Research shows number of samples need to be big enough to avoid excessive overfitting. See B section of the supplementary material of [23]. Note that the research is proven for AlexNet. While we don't have enough data to seek for the threshold which suffices to avoid overfitting for our choice of architecture, we assume for the ImageCLEF data sets [17], 20 per samples is not enough.

Data augmentation is widely used to reduce overfitting on image data [11]. This technique makes up the lower amount of samples and with some transformation increases the variation of the image. In our case, we've applied rotation, flipping, space shifting and channel shifting.

For our problemset, we experimented both to apply or to not apply data augmentation for the ImageCLEF dataset.

- 4) *Transfer Learning:* Since the output is different for our specific problem, we cannot apply the architecture of the pre-trained weights directly. At minimum we have to drop the softmax layer at the very end to retrain with our own dataset. The technique is known as Transfer Learning. The two well-known options for Transfer Learning are:

- 1) ConvNet as fixed feature extractor, and then classify with other Classifiers such as Logistic Regression/Softmax or SVM
- 2) Fine-tuning the ConvNet. For a N level structure, train the last H levels with the $N - H$ lower levels freeze (the

higher the level, the less overfitting to the target-datasets [20].

There are several options of architecture trained against the well-known ImageNet whose weights are available with keras framework.

- VGGNet, runner-up in ILSVRC 2014 (GoogLeNet is the winner of that year).
- ResNet, winner of ILSVRC 2015. The available weights are for ResNet50 with 50 layers (including Fully Connected layers) in total.

The ideality to choose an architecture of pre-trained weights is that it has been trained against original datasets that is similar to the target datasets. ImageNet has 14,197,122 images, 21841 synsets. Among them, there are 70 synsets that is related to leaves. Given the large number of syncsets, the weights trained against ImageNet is generalized enough that there will be relatively less bias. The weights are trained for the object image classification task of ImageNet, which is aligned with our task except that ours are more fine-grained leaf classification.

Further more, Jason et al. found that even features transferred from distant tasks are better than random weights [23]. Also note that ConvNet features are more generic in early layers and more original-dataset-specific in later layers[20].

After comparing preliminary results, we choose ResNet50 since ResNet50 gives better results and less overfitting. We believe this can be attributed to the fact that ResNet50 is deeper, but still having lower complexity[21]. It also generates lower dimension feature vector, which is likely due to the use of a more aggressive Average Pooling with a pool size of 7x7. This saves us from the effort to seek for reduction of dimensionality.

The ResNet is famous for it's deep layers[21], in our case, 50 layers, with 49 Conv layers and one FC layer on top.

Except for the first Conv layer, the rest 48 composes 16 “residual” blocks in 4 stages. The block within each stage has similar architecture, i.e. same input & output shape.

The possible approaches as forementioned, are either get the bottleneck features which is called CNN codes in the terms of transfer Learning, or freeze all the layers except for the last Residual Block and train with the target datasets. See Figure 5.

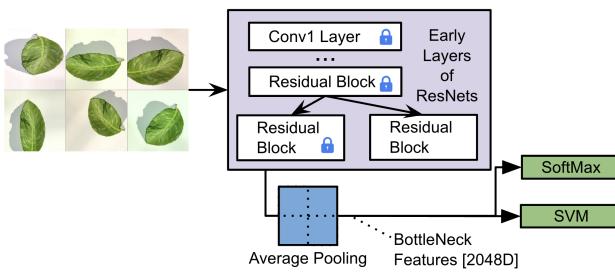


Fig. 5. Illustration of Transfer Learning with ResNet

Again, with the constrains of the time and computations, we choose to just extract the CNN codes.

We download the weights and load it into the preset model. And drop off the layer after the last Convolutional/Residual block. Feed the network with augmented data, the CNN codes (feature vector) we get from ResNet50 is 2048 dimensions.

With the help of visualization, intuitively we can see that each filter is seizing different features, roughly the outlines, the texture or the contour. The output from later layers is more abstract and global than the earlier layers. We can also see that some of the filters are completely dark, means for that particular filter, it doesn't response to certain feature of the image after rectification.

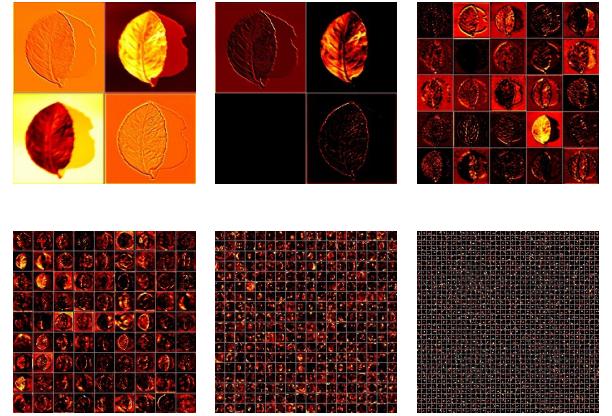


Fig. 6. Visualization of partial outputs for each stage of ResNet50. Left to right, top to bottom, they are from First Convolutional layer, the activation layer right after, Stage 2, Stage 3, Stage 4 and Stage 5 respectively. Note that the 4 outputs of Conv layer are corresponding to the 4 outputs of the activation layer

Thus, we can assume that the pre-trained weight is applicable to our problem, at least the early layers do generalize well to our feature space.

5) *Classification:* The feature vector is normalized with unit variance and we learn a linear SVM classifier (or others) from the features.

We applied grid search for the hyperparameters. For SVM, we searched for the C parameters with logarithmic scale .01, .1 and 1.

C. Discussion

How do you decide what type of transfer learning you should perform on a new dataset? This is a function of several factors, but the two most important ones are the size of the new dataset (small or big), and its similarity to the original dataset (e.g. ImageNet-like in terms of the content of images and the classes, or very different, such as microscope images). Keeping in mind that ConvNet features are more generic in early layers and more original-dataset-specific in later layers, here are some common rules of thumb for navigating the 4 major scenarios: [20]

- 1) similarity (transfer learning is good)
- 2) size for each class(fine tuning is not a good idea)

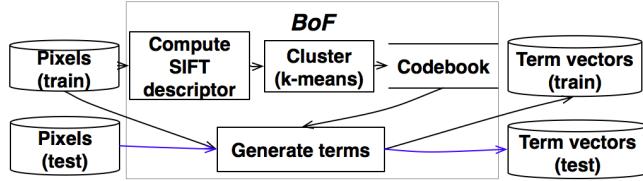


Fig. 7. Bag of Features

D. SIFT + Bag of Features (BoF)

Due to the simplicity and performance, this well established approach was taken at first. We prototyped the system with OpenCV libraries. Here is the illustration of the system [7]. Key points are densely sampled with a step size of 20 from the grays caned copies of the training set. Then we extract SIFT descriptor for each key point, which are clustered to build visual words via K-Means. Several codebook size (500, 1000, 3000) are used and we pick the one with the best validation results for each data set. To reduce computation complexity, we randomly select 100 training images to build the codebook. Finally, each train and test sample is represented with histograms of visual words, i.e. term vectors, and this is used as feature vectors for classification.

V. EXPERIMENTAL RESULTS

[Explain cross validation approach we used]

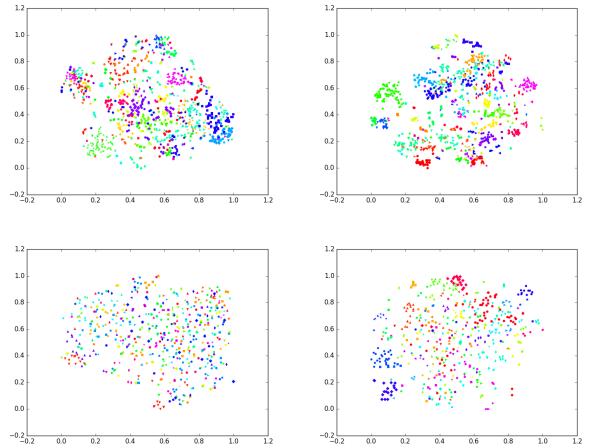
We use prediction rank-1 identification (i.e. accuracy) as our performance metric, which is defined as

$Accuracy = \frac{N_c}{N_t} \times 100\%$, where N_c represents the number of correct match and N_t is the total number of test samples.

[Test results as table.]

VI. DISCUSSIONS

- As expected, CNN codes off the shelf yields similar or better accuracy, compared to SIFT+BoF. Particularly traditional method suffers on noisy datasets. We use t-SNE [?] to visualize the feature vectors and it is clear that feature tend to spread further in Image CLEF Natural dataset. And the combination of Natural data set + traditional method gives the most sprawling representation in feature space.



Top left: Uniform BoF. Top right: Uniform Transfer Learning. Bottom left: Natural BoF. Bottom right: Natural Transfer Learning.

- To understand what stimulate neurons and how variations affect such stimulus, we visualize the feature map and deconv[28] of pool5 layer of ResNet for several pairs of data, using the visualization toolbox [?], in table I. Except for the last pair, all pairs come from the same species. And we can tell roughly that background clutter and color give quite different CNN codes.
 - Background clutter yields the biggest confusions for ConvNet. None of the top 5 ILSVRC labels match and the stimulated neurons are quite different.
 - Color and scale seem to activate slightly different sets of neurons.
- As an exercise for error analysis, we manually inspect the results of two preprocessing methods and pick the best out of these two for Image CLEF Uniform data set.

No preprocessing	K-means	Hand-picked
63.23%	58.79%	71.11%

 This clearly tells that the K-means preprocessing is not solid and the background noise contributes to the misclassification results.
- Looking at the confusion matrix, we believe the main causes for misclassification
 - Very fine differences between species, which is hard even for human experts
 - Noisy and possibly non-representative train data lead to overfitting,

VII. CONCLUSION AND FUTURE WORK

- Acquire more data and fine-tune ConvNet to solve overfitting problem
- Engage advanced techniques for image augmentation
- Explore state-of-art method to detect and locate leaf for Image CLEF natural leaf dataset.

REFERENCES

- S. Cho, D. Lee, and J. Jeong. Automation and emerging technologies: Weedplant discrimination by machine vision and artificial neural network. *Biosystems Engineering*, 83(3):275280, 2002.

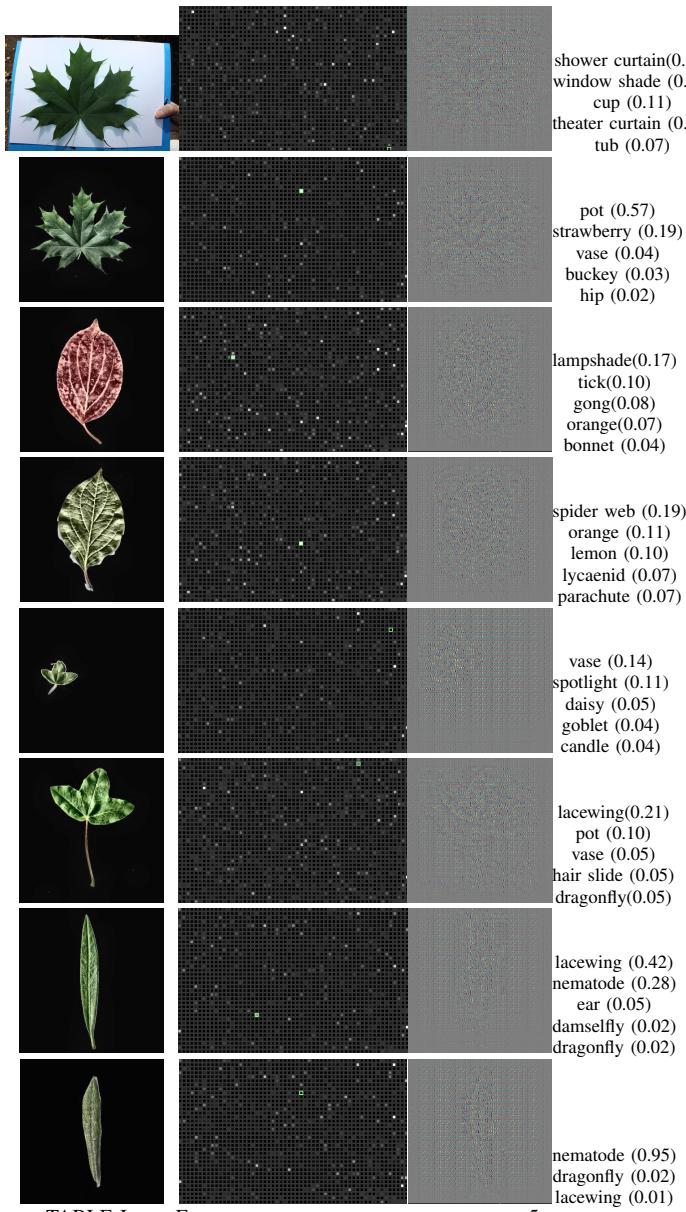


TABLE I. FEATURE MAP AND DECONV FOR POOL5 LAYER

- [2] Charles Mallah, James Cope, James Orwell. Plant Leaf Classification Using Probabilistic Integration of Shape, Texture and Margin Features. Signal Processing, Pattern Recognition and Applications, in press. 2013
- [3] Pedro F. B. Silva, Andre R.S. Marcal, Rubim M. Almeida da Silva. Evaluation of Features for Leaf Discrimination. 2013. Springer Lecture Notes in Computer Science, Vol. 7950, 197-204.
- [4] Itheri Yahiaoui, Nicolas Herve, and Nozha Boujemaa. Shape-based image retrieval in botanical collections, Lecture Notes in Computer Science including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics, vol. 4261 LNCS, pp 357-364, 2006.
- [5] Jassmann TJ, Tashakkori R, Parry RM (2015) Leaf classification utilizing a convolutional neural network. In: SoutheastCon
- [6] Neeraj Kumar, Peter N Belhumeur, Arijit Biswas, David W Jacobs, W John Kress, Ida C Lopez, and Joao VB Soares, Leafsnap: A computer vision system for automatic plant species identification, in ECCV, pp.