Introduction to image processing and computer vision

# Plant Species Recognition

Laboratory Project 2

Patryk Walczak

January 27, 2020

# Contents

# 1 Introduction

The main task of the project is an elaboration of a discriminative feature vector for the leaves of trees. Using computer vision processing find the vectors then learn the machine learning model with 80% of the data set and test the results on the remaining part. Six tree species each in the separated dictionary are stored in "isolated" dictionary. In each class is from 38 up to 97 images of leaves.



(a) Acer Circinatum
Vine Maplel

(b) Quercus Garryana
Oregon White Oak

(c) Acer Glabrum
Douglasii

(d) Quercus Kelloggii
California Black Oak

(e) Acer Macrophyllum
Big Leaf Maple

(f) Acer Negundo
Boxelder

Figure 1: Pictures of leaves

## 1.1  Project Description

For classification, I chose a random forest method, which is the ensemble learning method for classification, regression and other tasks that operates by constructing a multitude of decision trees. All images are analyzed by 6 descriptors, all of them are used to learn the algorithm and after that to recognise images from the test set. For all-purpose, I prepared 5 python files.

List of files :

- descriptors.py - contains all descriptors functions

- summary.py - uses to save all results in files

- separate_images.py - copy files to "train&test" directory and separate them to train and test sets

- train_test.py - the main file contains all train and test functions

- cross_validation.py - file contains function with cross-validation

# 2 Classification

In the section is described the process of classification leaves, using the descriptors and training and testing model.

## 2.1 Feature Descriptors

The subsection includes descriptions of all descriptors.

### 2.1.1 Hu Moments

Hu moments are used to calculate moments that are invariant to translation, scale, and rotation. The feature descriptor is in the first function in "descriptors.py". For clarify, an image moment is a certain particular weighted average (moment) of the image pixels' intensities or a function of such moments, usually chosen to have some attractive property or interpretation.

```
# Hu Moments
def fd_hu_moments(image):
gray = cv2.cvtColor(image, cv2.COLOR_BGR2GRAY)
feature = cv2.HuMoments(cv2.moments(gray)).flatten()
return feature
```

This is a simple function. It takes as an input the only image, changes the image's colors to black and white and computes Hu Moments. In the end, it flattens the results to one array.

### 2.1.2 Haralick Texture

Texture defines the consistency of patterns and colors in an image. Haralick suggested the use of grey level co-occurrence matrices (GLCM). The basic idea of GLCM - it looks for pairs of adjacent pixel values that occur in an image and keeps recording it over the entire image. The Haralick Texture is computed in the second function in "descriptors.py".

```
# Haralick Texture
def fd_haralick(image):
gray = cv2.cvtColor(image, cv2.COLOR_BGR2GRAY)
haralick = mahotas.features.haralick(gray).mean(axis=0)
return haralick
```

The function is simple, takes the only image as input, and return haralick descriptor. Firstly it changes the colors of the image to black and white. After that using mahotas library function computes the feature descriptor.

### 2.1.3 Color Histogram

In image processing and photography, a color histogram is a representation of the distribution of colors in an image. For digital images, a color histogram represents the number of pixels that have colors in each of a fixed list of color ranges, that span

the image's color space, the set of all possible colors. The ranges convert to the vector is the result of the third feature descriptor in "descriptors.py".

```
1 # Color Histogram
2 def fd_histogram(image):
3 hist = cv2.calcHist([image], [0, 1, 2], None, [8, 8, 8], [0, 256, 0,
      256, 0, 256])
4 cv2.normalize(hist, hist)
5 return hist.flatten()
```

The function is even easier than the previous two. The descriptor for computation needs the only image. Then even no changes of colours are required. Using the OpenCV library and the calcHist method the function compute the histogram. After that, the histogram is normalized and returned flattened.

### 2.1.4 Histogram of Oriented Gradients

The Histogram of Oriented Gradients descriptor technique counts occurrences of gradient orientation in localized portions of an image. The feature descriptor is computed as the fourth one in descriptors.py".

```
1 # Histogram of Oriented Gradients
2 def fd_histog(image):
3 gray = cv2.cvtColor(image, cv2.COLOR_BGR2GRAY)
4 h = feature.hog(gray, orientations=9, pixels_per_cell=(8, 8),
5 cells_per_block=(2, 2), transform_sqrt=True, block_norm="L1")
6 return h
```

As the rest descriptor functions, the one takes the only image as input also. Then changing the colors to black and white occurs. The feature library helps with the computing of the histogram of the oriented gradients. And the result of the hog method is returned.

### 2.1.5 Local Binary Pattern

In the whole project is only two local descriptors - the first one Local Binary Pattern, which computes a local representation of texture. The input image is changed to the black-and-white photo, then using feature library and its method local_binary_pattern compute the descriptor.

```
1 def fd_binary_pattern(image):
2 gray = cv2.cvtColor(image, cv2.COLOR_BGR2GRAY)
3 lbp = feature.local_binary_pattern(gray, 24, 8, method="uniform")
4 (hist, _) = np.histogram(lbp.ravel(),
5 bins=np.arange(0, 27),
6 range=(0, 26))
7 hist = hist.astype("float")
8 hist /= (hist.sum() + 1e-7)
9 return hist
```

### 2.1.6 Oriented FAST and rotated BRIEF

The second local descriptor is Oriented FAST and rotated BRIEF. For the one, the results were different for each kind of tree. To prevent the different length of the descriptor vector, after computing the descriptor the length is changed to 100. If the vector has less then 100 values the remaining is a sequence of zeros, in another case when vector has more than 100 values only the first hundred is taken as the result.

```python
# Oriented FAST and rotated BRIEF
def fd_orb(image):
orb = cv2.ORB_create(nfeatures=100)
keypoints = o.detect(image,None)
keypoints, descriptors = o.compute(image, keypoints)
descriptors=descriptors.flatten()
if len(descriptors)<100:
descriptors = np.zeros(100-len(descriptors),dtype=int)
else:
descriptors = np.array(descriptors[0:100])
return descriptors
```

### 2.1.7 Unused

I prepared Scale Invariant Feature Transform descriptor also but the method SIFT_create() from cv2.xfeatures2d returned the communication "This algorithm is patented and is excluded in this configuration". So that I decided not to use the descriptor in the project.

```python
# SIFT
def fd_SIFT(image):
gray = cv2.cvtColor(image, cv2.COLOR_BGR2GRAY)
freakExtractor = cv2.xfeatures2d.FREAK_create()
keypoints,descriptors= freakExtractor.compute(gray,None)
return keypoints
```

After searching for information, I found the free version of SIFT algorithm but it does not return any features. The algorithm throws an error because the machine learning model cannot learn from an empty feature vector. So as in the upper case, I decided to not use the descriptor in the project.

```python
# SIFT v2
def fd_SIFT_v2(image):
gray= cv2.cvtColor(image,cv2.COLOR_BGR2GRAY)
sift = cv2.xfeatures2d.SIFT_create()
kp = sift.detect(gray,None)
image=cv2.drawKeypoints(gray,kp,image)
cv2.show('sift_keypoints.jpg',image)
return kp
```

## 2.2 Cross validation

**Cross-validation is a process which predicts the ML model's accuracy.** In the project is used 10-fold cross-validation which means that the data set is divided into 10 subsets. The algorithm predicts the result 10 times, every time 9 sets are used to train and one for tests. Each iteration other set is the test one, so for each set is once used for test and nine times for train. The average of the ten results is returned.

```python
print("10-fold cross validation")
class_names = ["circinatum", "garryana", "glabrum", "kelloggii", "
    macrophyllum","negundo"]
all_descriptors = []
labels = []
for label in class_names:
for jpgfile in glob.glob("isolated/"+label+"/*.jpg"):
image = cv2.imread(jpgfile)
image = cv2.resize(image, (512,512))
```

Every image in the data set is loaded and resize to the same dimensions.

```python
fv_hu_moments  = fd_hu_moments(image)
fv_haralick    = fd_haralick(image)
fv_histogram   = fd_histogram(image)
fv_hog         = fd_histog(image)
fv_binpat      = fd_binary_pattern(image)
fv_orb         = fd_orb(image)
fv_SIFT        = fd_SIFT(image)
# save them in descriptors variable
 descriptors = np.hstack([fv_histogram , fv_haralick, fv_hu_moments
    , fv_hog , fv_binpat, fv_orb])
```

The stack of descriptors is saved in all_descriptors list and kind of tree is saved in labels list.

```python
# save labels and feature as the vectors
labels.append(label)
all_descriptors.append(descriptors)
```

Then from the library "sklearn.model_selection" take KFold function and prepare the algorithm for 10-fold cross-validation. After that function creates the ML model - Random Forest Classifier. In the end, the function computes the cross-validation and print mean and standard deviation of the ten results fo the 10-fold algorithm.

```python
# # 10-fold cross validation
kfold = KFold(n_splits=10)
rf = RandomForestClassifier(n_estimators=100, random_state=9)
cv_results = cross_val_score(rf, train_data, train_labels, cv=kfold,
    scoring="accuracy")
msg = "%s: %f (%f)" % ('RF', cv_results.mean(), cv_results.std())
print(msg)
```

## 2.3 Division of data set

For classification purposes, the data set has to be split into two parts. The first with 80% of all images are used to train the method. The second with remaining 20% is used for test of the method.

In "separete_images.py" I prepared to function for split the data set.

```
class_names = ["circinatum", "garryana", "glabrum", "kelloggii", "
    macrophyllum","negundo"]
directory = "train&test"

def separeteImages():
```

Firstly, the program creates new empty directories to store the images.

```
# create new empty directories
if os.path.exists(directory):
shutil.rmtree(directory, ignore_errors=True)
os.makedirs(directory)
```

Next, for every kind of tree create test and train directories.

```
# for every spiece
for image in class_names:
# prepare two directories
# one for train images
train_dir = directory+"/"+image
os.mkdir(train_dir)
train_dir += "/train"
os.mkdir(train_dir)
# one for test images
test_dir = directory+"/"+image+"/test"
os.mkdir(test_dir)
```

After that, the program start splitting all photos in the set

```
# for every image in data set
i=0
for jpgfile in glob.glob("isolated/"+image+"/*.jpg"):
if(i%5==0):
shutil.copy(jpgfile, test_dir)
else:
shutil.copy(jpgfile, train_dir)
i+=1
```

Every image with index number equal modulo to 5 is saved in the test directory, and the rest are stored in train directory. For comparison of the descriptors, I preformed the code only once at the start of the project. The fact provides that for all results algorithm always has the same input data.

## 2.4 Train and test

**The section covers the most important part of the project. Here is the description of the main algorithm which learn the ML model and after that tests the results. The algorithm is included in "train_test.py"**

At the start, the program reads all imagines in train directories from all species and resize all of them to common dimensions.

```
1  all_descriptors = []
2  labels = []
3  # learning from the train directories
4  for label in class_names:
5  for file in glob.glob('train&test/'+(label)+'/train/*.jpg'):
6  # read the image and resize it to the common size
7  image = cv2.imread(file)
8  image = cv2.resize(image, (512,512))
```

For each image, all descriptors are computed and saved in all_descriptors list, while the species names of the images are stored in labels list.

```
1  # compute descriptors
2  fv_hu_moments = fd_hu_moments(image)
3  fv_haralick   = fd_haralick(image)
4  fv_histogram  = fd_histogram(image)
5  fv_hog        = fd_histog(image)
6  fv_binpat     = fd_binary_pattern(image)
7  fv_orb        = fd_orb(image)
8  # fv_SIFT      = fd_SIFT(image)
9  # save them in descriptors variable
10 descriptors = np.hstack([fv_histogram, fv_haralick, fv_hu_moments,
       fv_hog, fv_binpat, fv_orb])
11 # save labels and feature as the vectors
12 labels.append(label)
13 all_descriptors.append(descriptors)
```

When the training set of data completed, the program creates the Random Forest Classifier and learns the ML model.

```
1  # Random Forests Classifier
2  clf = RandomForestClassifier(n_estimators=100, random_state=9)
3  # learning the model
4  clf.fit(all_descriptors, labels)
```

When learning is done, the algorithm starts testing section by reads all images and resize them.

```
1  for label in class_names:
2  for file in glob.glob('train&test/'+(label)+'/test/*.jpg'):
3  # read the image and resize it to the common size
4  image = cv2.imread(file)
5  image = cv2.resize(image, (512,512))
```

Once again program computes all descriptors but in the case for the testing images.

```
1          # compute descriptors
2          fv_hu_moments = fd_hu_moments(image)
```

```
3        fv_haralick    = fd_haralick ( image )
4        fv_histogram   = fd_histogram ( image )
5        fd_hog         = fd_histog ( image )
6        fd_binpat      = fd_binary_pattern ( image )
7        fv_orb         = fd_orb ( image )
8        # fv_SIFT        = fd_SIFT ( image )
9        # save them in descriptors variable
10       descriptors = np . hstack ([ fv_histogram , fv_haralick ,
    fv_hu_moments , fv_hog , fv_binpat , fv_orb ])
```

Next step is the prediction and checking the results. Name of the kind of tree is printed out with the predicted name. User can easily and in real-time check the results.

```
1  prediction = clf . predict ( descriptors . reshape (1 , -1)) [0]
2  print ( label , prediction )
3  results . append (( label , prediction ))
```

All of the results are stored in the results list. The reason for that is an optional function which can save the data to files and computes the mean. The function "save_and_summary" is described in the next subsection.

```
1  save_and_summary ( results )
```

## 2.5 Summary

**For saving results for each image, each kind of tree and for all data set I prepared summary.py file.**

At first, the program creates empty files in directory "test&file". AllResults.txt file includes all results in form "kind of tree" → "prediction", and Summary.txt file includes a summary for each class and whole data set.

```
1  # function for save results and summary
2  def save_and_summary ( results ):
3  # create files to save results
4  if not os.path.exists ("./train&test"):
5  os.makedirs ("./train&test")
6  AllResults = open ("./train&test/AllResults.txt" , 'w')
7  Summary = open ("./train&test/Summary.txt" , 'w')
```

```
Summary
circinatum average is equal 0.9285714285714286
garryana average is equal 1.0
glabrum average is equal 1.0
kelloggii average is equal 1.0
macrophyllum average is equal 1.0
negundo average is equal 1.0
Average for whole dataset is equal 0.989010989010989
```

```
macrophyllum -> macrophyllum
macrophyllum -> macrophyllum
macrophyllum -> macrophyllum
macrophyllum -> macrophyllum
negundo -> negundo
negundo -> negundo
negundo -> negundo
negundo -> negundo
negundo -> negundo
```

Variable "results" is a list of tuples with the correct answer at the 1st position and prediction of the model at the 2nd position in each tuple.

```
1  # from train&test.py
2  results.append ((label,prediction))
```

Then create two dictionaries and save the results to AllResults.txt. In the same loop compute the number of correct predictions.

```
 1  # create 2 dictionaries for computing correct preditions
 2  num = dict()
 3  den = dict()
 4  for name in class_names:
 5  num[name] = 0
 6  den[name] = 0
 7  for result in results:
 8  # write all preditions to AllResults file
 9  AllResults.write(' -> '.join(result)+'\n')
10  # compute correctness for each spiecie
11  den[result[0]] += 1
12  if result[0] == result[1]:
13  num[result[0]] += 1
```

In the next step, the program starts summary. Firstly, compute the result for each kind of tree.

```
 1  n = 0
 2  d = 0
 3  # Summary
 4  Summary.write("Summary\n")
 5  for name in class_names:
 6  # write results for each class
 7  if den[name] != 0:
 8  line = name+" average is equal "+ str(num[name]/den[name])
 9  Summary.write(line+"\n")
10  print(line)
11  n += num[name]
12  d += den[name]
```

Then it computes mean for all data set.

```
 1  if d != 0:
 2  # write result for whole data set
 3  line = "Average for whole dataset is equal "+ str(n/d)
 4  Summary.write(line+"\n")
 5  print(line)
```

At the end close the files.

# 3 Results

## 3.1 10-fold cross validataion

**For better results I prepared 10-fold cross validataion for whole set of data.**

The 10-fold algorithm divides the set into 10 parts and checks the results of machine learning for each set as the test one and rest as the training. Then mean of all 10 outputs it the finite result. The image below shows the results for all descriptors.

```
10-fold cross validation
Computed descriptors of circinatum
Computed descriptors of garryana
Computed descriptors of glabrum
Computed descriptors of kelloggii
Computed descriptors of macrophyllum
Computed descriptors of negundo
RF: 0.391919 (0.292805)
```

The mean is only 40%. However, for the 3 descriptors (Hu moments, Haralick Texture, Color Histogram) the cross-validation result has the maximum mean which is equal to 89%.

```
10-fold cross validation
Computed descriptors of circinatum
Computed descriptors of garryana
Computed descriptors of glabrum
Computed descriptors of kelloggii
Computed descriptors of macrophyllum
Computed descriptors of negundo
RF: 0.886515 (0.258958)
```

The most interesting fact is that if I had used only Color Histogram descriptor the results were close to the highest result.

```
10-fold cross validation
Computed descriptors of circinatum
Computed descriptors of garryana
Computed descriptors of glabrum
Computed descriptors of kelloggii
Computed descriptors of macrophyllum
Computed descriptors of negundo
RF: 0.884242 (0.258068)
```

So after adding 3 global and 2 local descriptors the result increases by only 0,2%. The 5 remaining descriptors have result between 40% and 70%. The results are around 90% of accuracy as the mean of 10-fold cross-validation.

## 3.2  80% training / 20% testing

**In the project description stayed the statement that the whole data set should be divided into two parts. First with 80% of the images as the training and rest 20% as the test set. Using "separeteImages()" function which is described in the 2.3 section I divide the data and perform the test.**

For the results, I prepared "save_and_summary( results )" function which is described in the 2.5 section. All the results are saved as text files. Moreover, each species has the mean results and each prediction is saved also. In the report, I show only the summary of the few sets of descriptors. Firstly the result for all descriptors.

```
circinatum average is equal 0.0
garryana average is equal 0.0
glabrum average is equal 0.0
kelloggii average is equal 0.0
macrophyllum average is equal 0.0
negundo average is equal 1.0
Average for whole dataset is equal 0.08791208791208792
```

The result is low - only 9%. For the negundo kind, the result is equal to 100% but this is the only class without 0% as a result. However below is shown the image with the result using only 3 descriptors (Hu moments, Haralick Texture, Color Histogram).

```
circinatum average is equal 0.9285714285714286
garryana average is equal 1.0
glabrum average is equal 1.0
kelloggii average is equal 1.0
macrophyllum average is equal 1.0
negundo average is equal 1.0
Average for whole dataset is equal 0.989010989010989
```

For the 3 descriptors, the results are much better - around 99%. And only for the circinatum kind of tree the results are not 100% correct. And one more result of the comparison is presented below.

```
circinatum average is equal 0.9285714285714286
garryana average is equal 1.0
glabrum average is equal 1.0
kelloggii average is equal 1.0
macrophyllum average is equal 1.0
negundo average is equal 1.0
Average for whole dataset is equal 0.989010989010989
```

The results look like 2nd time the same image but this one is different - only Color Histogram was used to obtain the results. So in the case, the 5 descriptors (3 globals and 2 locals) do not affect the best result. The machine learning model correctly interprets the same number of images using only Color Histogram, the same number using the 3 descriptors and obtain better result then all set of descriptors. The 10-fold cross-validation shows the small difference but for the 80% training / 20% testing the results are the same.

# 4 Source Code

## 4.1 descriptors.py

```python
1  import cv2
2  import mahotas
3  import numpy as np
4  from skimage import feature
5
6  # Hu Moments
7  def fd_hu_moments(image):
8      gray = cv2.cvtColor(image, cv2.COLOR_BGR2GRAY)
9      feature = cv2.HuMoments(cv2.moments(gray)).flatten()
10     return feature
11
12 # Haralick Texture
13 def fd_haralick(image):
14     gray = cv2.cvtColor(image, cv2.COLOR_BGR2GRAY)
15     haralick = mahotas.features.haralick(gray).mean(axis=0)
16     return haralick
17
18 # Color Histogram
19 def fd_histogram(image):
20     hist  = cv2.calcHist([image], [0, 1, 2], None, [8, 8, 8], [0,
    256, 0, 256, 0, 256])
21     cv2.normalize(hist, hist)
22     return hist.flatten()
23
24 # Histogram of Oriented Gradients
25 def fd_histog(image):
26     gray = cv2.cvtColor(image, cv2.COLOR_BGR2GRAY)
27     h = feature.hog(gray, orientations=9, pixels_per_cell=(8, 8),
28             cells_per_block=(2, 2), transform_sqrt=True, block_norm=
    "L1")
29     return h
30
31 # Local Binary Pattern
32 def fd_binary_pattern(image):
33     gray = cv2.cvtColor(image, cv2.COLOR_BGR2GRAY)
34     lbp = feature.local_binary_pattern(gray, 24, 8, method="uniform"
    )
35     (hist, _) = np.histogram(lbp.ravel(),
36         bins=np.arange(0, 27),
37         range=(0, 26))
38     hist = hist.astype("float")
39     hist /= (hist.sum() + 1e-7)
40     return hist
41
42 # Oriented FAST and rotated BRIEF
43 def fd_orb(image):
44     orb = cv2.ORB_create(nfeatures=100)
45     keypoints = orb.detect(image,None)
46     keypoints, descriptors = orb.compute(image, keypoints)
47     descriptors=descriptors.flatten()
```

```
48    if len(descriptors)<100:
49        descriptors = np.zeros(100-len(descriptors),dtype=int)
50    else:
51        descriptors = np.array(descriptors[0:100])
52    return descriptors
53
54 # SIFT
55 def fd_SIFT(image):
56    gray = cv2.cvtColor(image, cv2.COLOR_BGR2GRAY)
57    freakExtractor = cv2.xfeatures2d.FREAK_create()
58    keypoints,descriptors= freakExtractor.compute(gray,None)
59    return keypoints
60
61 # SIFT v2
62 def fd_SIFT_v2(image):
63    gray= cv2.cvtColor(image,cv2.COLOR_BGR2GRAY)
64    sift = cv2.xfeatures2d.SIFT_create()
65    kp = sift.detect(gray,None)
66    image=cv2.drawKeypoints(gray,kp,image)
67    cv2.show('sift_keypoints.jpg',image)
68    return kp
69
70 # for tests
71 if __name__ == "__main__":
72    image = cv2.imread('.\isolated/negundo\l20.jpg')
73    print(fd_histog(image))
```

## 4.2 summary.py

```
1 import os
2
3 # list of all spiecies
4 class_names = ["circinatum", "garryana", "glabrum", "kelloggii", "
     macrophyllum","negundo"]
5
6 # function for save results and summary
7 def save_and_summary( results ):
8    # create files to save results
9    if not os.path.exists("./train&test"):
10        os.makedirs("./train&test")
11    AllResults = open("./train&test/AllResults.txt" , 'w')
12    Summary = open("./train&test/Summary.txt" , 'w')
13    # create 2 dictionaries for computing correct preditions
14    num = dict()
15    den = dict()
16    for name in class_names:
17        num[name] = 0
18        den[name] = 0
19    for result in results:
20        # write all preditions to AllResults file
21        AllResults.write(' -> '.join(result)+'\n')
22        # compute correctness for each spiecie
23        den[result[0]] += 1
24        if result[0] == result[1]:
25            num[result[0]] += 1
```

```
26       n = 0
27       d = 0
28       # Summary
29       Summary.write("Summary\n")
30       for name in class_names:
31           # write results for each class
32           if den[name] != 0:
33               line = name+" average is equal "+ str(num[name]/den[name
     ])
34               Summary.write(line+"\n")
35               print(line)
36           n += num[name]
37           d += den[name]
38       if d != 0:
39           # write result for whole data set
40           line = "Average for whole dataset is equal "+ str(n/d)
41           Summary.write(line+"\n")
42           print(line)
43       # close files
44       AllResults.close()
45       Summary.close()
46
47 if __name__ == "__main__":
48     results = list()
49     save_and_summary(results)
```

## 4.3   separate_images.py

```
1 import os
2 import glob
3 import shutil
4
5 class_names = ["circinatum", "garryana", "glabrum", "kelloggii", "
     macrophyllum","negundo"]
6 directory = "train&test"
7
8 def separeteImages():
9     # create new empty directories
10    if os.path.exists(directory):
11        shutil.rmtree(directory, ignore_errors=True)
12    os.makedirs(directory)
13    # for every spiece
14    for image in class_names:
15        # prepare two directories
16        # one for train images
17        train_dir = directory+"/"+image
18        os.mkdir(train_dir)
19        train_dir += "/train"
20        os.mkdir(train_dir)
21        # one for test images
22        test_dir = directory+"/"+image+"/test"
23        os.mkdir(test_dir)
24        # for every image in data set
25        i=0
26        for jpgfile in glob.glob("isolated/"+image+"/*.jpg"):
```

```
27            if(i%5==0):
28                 shutil.copy(jpgfile, test_dir)
29            else:
30                 shutil.copy(jpgfile, train_dir)
31            i+=1
32
33 if __name__ == "__main__":
34     separeteImages()
```

## 4.4   train_test.py

```
1  import glob
2  import cv2
3  import numpy as np
4  from sklearn.ensemble import RandomForestClassifier
5  from sklearn.model_selection import KFold, cross_val_score
6  from descriptors import fd_hu_moments
7  from descriptors import fd_haralick
8  from descriptors import fd_histogram
9  from descriptors import fd_histog
10 from descriptors import fd_binary_pattern
11 from descriptors import fd_SIFT
12 from descriptors import fd_orb
13 from summary import save_and_summary
14
15 class_names = ["circinatum", "garryana", "glabrum", "kelloggii", "
      macrophyllum","negundo"]
16
17 all_descriptors = []
18 labels          = []
19 # learning from the train directories
20 for label in class_names:
21     for file in glob.glob('train&test/'+(label)+'/train/*.jpg'):
22         # read the image and resize it to the common size
23         image = cv2.imread(file)
24         image = cv2.resize(image, (512,512))
25         # compute descriptors
26         fv_hu_moments = fd_hu_moments(image)
27         fv_haralick   = fd_haralick(image)
28         fv_histogram  = fd_histogram(image)
29         fv_hog        = fd_histog(image)
30         fv_binpat     = fd_binary_pattern(image)
31         fv_orb        = fd_orb(image)
32         # fv_SIFT       = fd_SIFT(image)
33         # save them in descriptors variable
34         descriptors = np.hstack([fv_histogram, fv_haralick,
      fv_hu_moments, fv_hog, fv_binpat, fv_orb])
35         # save labels and feature as the vectors
36         labels.append(label)
37         all_descriptors.append(descriptors)
38     print("Computed descriptors of "+label)
39
40 # Random Forests Classifier
41 clf  = RandomForestClassifier(n_estimators=100, random_state=9)
42 # learning the model
```

```
43  clf.fit(all_descriptors, labels)
44
45  results = list()
46  print("\nPredictions\n")
47  for label in class_names:
48      for file in glob.glob('train&test/'+(label)+'/test/*.jpg'):
49          # read the image and resize it to the common size
50          image = cv2.imread(file)
51          image = cv2.resize(image, (512,512))
52          # compute descriptors
53          fv_hu_moments = fd_hu_moments(image)
54          fv_haralick   = fd_haralick(image)
55          fv_histogram  = fd_histogram(image)
56          fd_hog        = fd_histog(image)
57          fd_binpat     = fd_binary_pattern(image)
58          fv_orb        = fd_orb(image)
59          # fv_SIFT      = fd_SIFT(image)
60          # save them in descriptors variable
61          descriptors = np.hstack([fv_histogram, fv_haralick,
        fv_hu_moments, fv_hog, fv_binpat, fv_orb])
62          prediction = clf.predict(descriptors.reshape(1,-1))[0]
63          print(label,prediction)
64          results.append((label,prediction))
65
66  save_and_summary(results)
```

## 4.5 cross_validation.py

```
1   from descriptors import fd_hu_moments
2   from descriptors import fd_haralick
3   from descriptors import fd_histogram
4   from descriptors import fd_histog
5   from descriptors import fd_binary_pattern
6   from descriptors import fd_SIFT
7   from descriptors import fd_orb
8   import glob
9   import cv2
10  import numpy as np
11  from sklearn.model_selection import train_test_split,
        cross_val_score
12  from sklearn.model_selection import KFold
13  from sklearn.ensemble import RandomForestClassifier
14
15  # 10-fold cross validation
16  def cross_validation():
17      print("10-fold cross validation")
18      class_names = ["circinatum", "garryana", "glabrum", "kelloggii",
        "macrophyllum","negundo"]
19      all_descriptors = []
20      labels          = []
21      for label in class_names:
22          for jpgfile in glob.glob("isolated/"+label+"/*.jpg"):
23              image = cv2.imread(jpgfile)
24              image = cv2.resize(image, (512,512))
25              # compute descriptors
```

```
26              fv_hu_moments = fd_hu_moments(image)
27              fv_haralick   = fd_haralick(image)
28              fv_histogram  = fd_histogram(image)
29              fv_hog        = fd_histog(image)
30              fv_binpat     = fd_binary_pattern(image)
31              fv_orb        = fd_orb(image)
32              fv_SIFT       = fd_SIFT(image)
33              # save them in descriptors variable
34              descriptors = np.hstack([fv_histogram , fv_haralick ,
      fv_hu_moments , fv_hog , fv_binpat, fv_orb])
35              # save labels and feature as the vectors
36              labels.append(label)
37              all_descriptors.append(descriptors)
38          print("Computed descriptors of "+ label)
39
40      # # 10-fold cross validation
41      kfold = KFold(n_splits=10)
42      rf = RandomForestClassifier(n_estimators=100, random_state=9)
43      cv_results = cross_val_score(rf, all_descriptors , labels , cv=
      kfold , scoring="accuracy")
44      msg = "%s: %f (%f)" % ('RF', cv_results.mean(), cv_results.std()
      )
45      print(msg)
46
47  if __name__ == "__main__":
48      cross_validation()
```

# 5 References

1. https://en.wikipedia.org/wiki/Random_forest

2. https://gogul.dev/software/image-classification-python

3. https://www.learnopencv.com/histogram-of-oriented-gradients/

4. https://opencv-python-tutroals.readthedocs.io/en/latest/py_tutorials/py_feature2d/py_orb/py_orb.html

5. https://en.wikipedia.org/wiki/Image_moment

6. http://wiki.awf.forst.uni-goettingen.de/wiki/index.php/Haralick_Texture

7. https://software.intel.com/en-us/ipp-dev-reference-histogram-of-oriented-\
gradients-hog-descriptor

8. https://en.wikipedia.org/wiki/Color_histogram

9. https://software.intel.com/en-us/ipp-dev-reference-histogram-of-oriented-\
gradients-hog-descriptor

10. https://www.pyimagesearch.com/2015/12/07/local-binary-patterns-with-python-\
opencv/