

Sports Celebrity Image Classification

Special thanks to Debjyoti Paul (My data scientist friend at Amazon) for help with this project

```
import numpy as np
import cv2
import matplotlib
from matplotlib import pyplot as plt
%matplotlib inline
```

(1) Preprocessing: Detect face and eyes

When we look at any image, most of the time we identify a person using a face. An image might contain multiple faces, also the face can be obstructed and not clear. The first step in our pre-processing pipeline is to detect faces from an image. Once face is detected, we will detect eyes, if two eyes are detected then only we keep that image otherwise discard it.

Now how do you detect face and eyes?

We will use haar cascade from opencv for this. Here is an article on this: https://opencv-python-tutroals.readthedocs.io/en/latest/py_tutorials/py_objdetect/py_face_detection/py_face_detection.html?highlight=haar

```
img = cv2.imread('./test_images/sharapova1.jpg')
img.shape
```

(555, 700, 3)

```
plt.imshow(img)
```

<matplotlib.image.AxesImage at 0x25581716b00>



```
gray = cv2.cvtColor(img, cv2.COLOR_BGR2GRAY)
gray.shape
```

```
(555, 700)
```

```
gray
```

```
array([[175, 175, 175, ..., 176, 175, 174],
       [175, 175, 175, ..., 177, 175, 174],
       [175, 175, 175, ..., 177, 176, 174],
       ...,
       [ 84,  87,  88, ..., 113, 113, 113],
       [ 88,  89,  90, ..., 113, 113, 113],
       [ 93,  91,  91, ..., 112, 112, 112]], dtype=uint8)
```

```
plt.imshow(gray, cmap='gray')
```

```
<matplotlib.image.AxesImage at 0x255817c0588>
```



```
face_cascade = cv2.CascadeClassifier('./opencv/haarcascades/haarcascade_frontalface_default.xml')
eye_cascade = cv2.CascadeClassifier('./opencv/haarcascades/haarcascade_eye.xml')

faces = face_cascade.detectMultiScale(gray, 1.3, 5)
faces
```

```
array([[352, 38, 233, 233]], dtype=int32)
```

```
(x,y,w,h) = faces[0]
x,y,w,h
```

```
(352, 38, 233, 233)
```

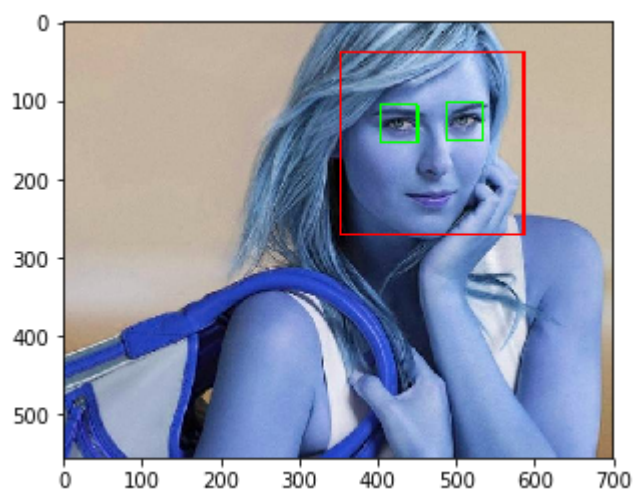
```
face_img = cv2.rectangle(img,(x,y),(x+w,y+h),(255,0,0),2)
plt.imshow(face_img)
```

```
<matplotlib.image.AxesImage at 0x25583bc41d0>
```



```
cv2.destroyAllWindows()
for (x,y,w,h) in faces:
    face_img = cv2.rectangle(img,(x,y),(x+w,y+h),(255,0,0),2)
    roi_gray = gray[y:y+h, x:x+w]
    roi_color = face_img[y:y+h, x:x+w]
    eyes = eye_cascade.detectMultiScale(roi_gray)
    for (ex,ey,ew,eh) in eyes:
        cv2.rectangle(roi_color,(ex,ey),(ex+ew,ey+eh),(0,255,0),2)

plt.figure()
plt.imshow(face_img, cmap='gray')
plt.show()
```



(2) Preprocessing: Crop the facial region of the image

```
%matplotlib inline
plt.imshow(roi_color, cmap='gray')
```

<matplotlib.image.AxesImage at 0x255855c17f0>



```
cropped_img = np.array(roi_color)
cropped_img.shape
```

```
(233, 233, 3)
```

(3) Preprocessing: Use wavelet transform as a feature for training our model

In wavelet transformed image, you can see edges clearly and that can give us clues on various facial features such as eyes, nose, lips etc

Wavelet transform

```
import numpy as np
import pywt
import cv2

def w2d(img, mode='haar', level=1):
    imArray = img
    #Datatype conversions
    #convert to grayscale
    imArray = cv2.cvtColor( imArray,cv2.COLOR_RGB2GRAY )
    #convert to float
    imArray = np.float32(imArray)
    imArray /= 255;
    # compute coefficients
    coeffs=pywt.wavedec2(imArray, mode, level=level)

    #Process Coefficients
    coeffs_H=list(coeffs)
    coeffs_H[0] *= 0;

    # reconstruction
    imArray_H=pywt.waverec2(coeffs_H, mode);
    imArray_H *= 255;
    imArray_H = np.uint8(imArray_H)

    return imArray_H
```

```
im_har = w2d(cropped_img,'db1',5)
plt.imshow(im_har, cmap='gray')
```

```
<matplotlib.image.AxesImage at 0x2558858db00>
```



You can see above a wavelet transformed image that gives clues on facial features such as eyes, nose, lips etc. This along with raw pixel image can be used as an input for our classifier

(3) Preprocessing: Load image, detect face. If eyes ≥ 2 , then save and crop the face region

Lets write a python function that can take input image and returns cropped image (if face and eyes ≥ 2 are detected)

```
def get_cropped_image_if_2_eyes(image_path):
    img = cv2.imread(image_path)
    gray = cv2.cvtColor(img, cv2.COLOR_BGR2GRAY)
    faces = face_cascade.detectMultiScale(gray, 1.3, 5)
    for (x,y,w,h) in faces:
        roi_gray = gray[y:y+h, x:x+w]
        roi_color = img[y:y+h, x:x+w]
        eyes = eye_cascade.detectMultiScale(roi_gray)
        if len(eyes) >= 2:
            return roi_color
```

```
original_image = cv2.imread('./test_images/sharapova1.jpg')
plt.imshow(original_image)
```

<matplotlib.image.AxesImage at 0x255885ef208>



```
cropped_image = get_cropped_image_if_2_eyes('./test_images/sharapova1.jpg')
plt.imshow(cropped_image)
```

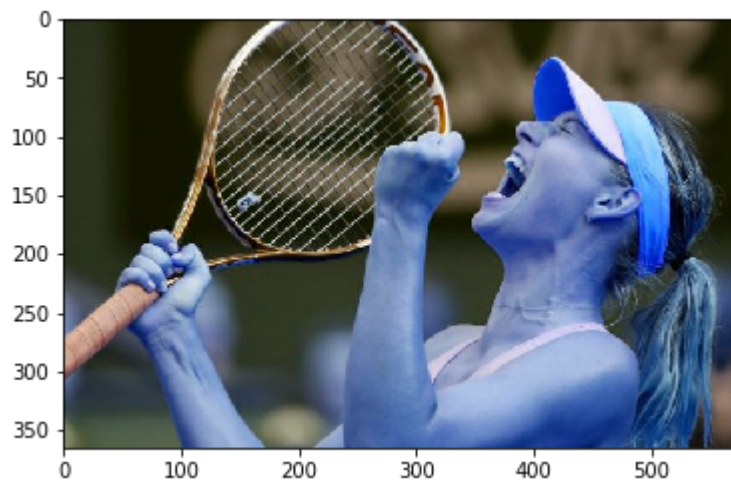
<matplotlib.image.AxesImage at 0x255888ab668>



In below image face is not very clear and it doesn't have two eyes clearly visible

```
org_image_obstructed = cv2.imread('./test_images/sharapova2.jpg')
plt.imshow(org_image_obstructed)
```

<matplotlib.image.AxesImage at 0x255888fcbe0>



```
cropped_image_no_2_eyes = get_cropped_image_if_2_eyes('./test_images/sharapova2.j
cropped_image_no_2_eyes
```

Above `cropped_image_no_2_eyes` is `None` which means we should ignore this image and we will not use such image for model training

```
path_to_data = "./dataset/"
path_to_cr_data = "./dataset/cropped/"
```

```
import os
img_dirs = []
for entry in os.scandir(path_to_data):
    if entry.is_dir():
        img_dirs.append(entry.path)
```

```
img_dirs
```

```
[ './dataset/cropped',
  './dataset/lionel_messi',
  './dataset/maria_sharapova',
  './dataset/roger_federer',
  './dataset/serena_williams',
  './dataset/virat_kohli']
```

Go through all images in dataset folder and create cropped images for them. There will be cropped folder inside dataset folder after you run this code

```
import shutil
if os.path.exists(path_to_cr_data):
    shutil.rmtree(path_to_cr_data)
os.mkdir(path_to_cr_data)
```

```
cropped_image_dirs = []
celebrity_file_names_dict = {}
for img_dir in img_dirs:
    count = 1
    celebrity_name = img_dir.split('/')[-1]
    celebrity_file_names_dict[celebrity_name] = []
    for entry in os.scandir(img_dir):
        roi_color = get_cropped_image_if_2_eyes(entry.path)
        if roi_color is not None:
            cropped_folder = path_to_cr_data + celebrity_name
            if not os.path.exists(cropped_folder):
                os.makedirs(cropped_folder)
            cropped_image_dirs.append(cropped_folder)
            print("Generating cropped images in folder: ",cropped_folder)
            cropped_file_name = celebrity_name + str(count) + ".png"
            cropped_file_path = cropped_folder + "/" + cropped_file_name
            cv2.imwrite(cropped_file_path, roi_color)
            celebrity_file_names_dict[celebrity_name].append(cropped_file_path)
            count += 1
```

```
Generating cropped images in folder: ./dataset/cropped/lionel_messi
Generating cropped images in folder: ./dataset/cropped/maria_sharapova
Generating cropped images in folder: ./dataset/cropped/roger_federer
Generating cropped images in folder: ./dataset/cropped/serena_williams
Generating cropped images in folder: ./dataset/cropped/virat_kohli
```

Now you should have cropped folder under datasets folder that contains cropped images

Manually examine cropped folder and delete any unwanted images

```
celebrity_file_names_dict = {}
for img_dir in cropped_image_dirs:
    celebrity_name = img_dir.split('/')[-1]
    file_list = []
    for entry in os.scandir(img_dir):
        file_list.append(entry.path)
    celebrity_file_names_dict[celebrity_name] = file_list
celebrity_file_names_dict
```

```
{'lionel_messi': ['./dataset/cropped/lionel_messi\\lionel_messi1.png',
 './dataset/cropped/lionel_messi\\lionel_messi10.png',
 './dataset/cropped/lionel_messi\\lionel_messi11.png',
 './dataset/cropped/lionel_messi\\lionel_messi13.png',
 './dataset/cropped/lionel_messi\\lionel_messi14.png',
 './dataset/cropped/lionel_messi\\lionel_messi15.png',
 './dataset/cropped/lionel_messi\\lionel_messi16.png',
 './dataset/cropped/lionel_messi\\lionel_messi17.png',
```


[illegible]

[illegible]

```
'./dataset/cropped/virat_kohli\\virat_kohli16.png',  
'./dataset/cropped/virat_kohli\\virat_kohli17.png',  
'./dataset/cropped/virat_kohli\\virat_kohli18.png',  
'./dataset/cropped/virat_kohli\\virat_kohli19.png',  
'./dataset/cropped/virat_kohli\\virat_kohli2.png',  
'./dataset/cropped/virat_kohli\\virat_kohli20.png',  
'./dataset/cropped/virat_kohli\\virat_kohli21.png',  
'./dataset/cropped/virat_kohli\\virat_kohli23.png',  
'./dataset/cropped/virat_kohli\\virat_kohli25.png',  
'./dataset/cropped/virat_kohli\\virat_kohli26.png',  
'./dataset/cropped/virat_kohli\\virat_kohli27.png',  
'./dataset/cropped/virat_kohli\\virat_kohli28.png',  
'./dataset/cropped/virat_kohli\\virat_kohli30.png',  
'./dataset/cropped/virat_kohli\\virat_kohli31.png',  
'./dataset/cropped/virat_kohli\\virat_kohli32.png',  
'./dataset/cropped/virat_kohli\\virat_kohli33.png',  
'./dataset/cropped/virat_kohli\\virat_kohli34.png',  
'./dataset/cropped/virat_kohli\\virat_kohli35.png',  
'./dataset/cropped/virat_kohli\\virat_kohli36.png',  
'./dataset/cropped/virat_kohli\\virat_kohli37.png',  
'./dataset/cropped/virat_kohli\\virat_kohli38.png',  
'./dataset/cropped/virat_kohli\\virat_kohli4.png',  
'./dataset/cropped/virat_kohli\\virat_kohli40.png',  
'./dataset/cropped/virat_kohli\\virat_kohli41.png',  
'./dataset/cropped/virat_kohli\\virat_kohli42.png',  
'./dataset/cropped/virat_kohli\\virat_kohli44.png',  
'./dataset/cropped/virat_kohli\\virat_kohli45.png',  
'./dataset/cropped/virat_kohli\\virat_kohli47.png',  
'./dataset/cropped/virat_kohli\\virat_kohli48.png',  
'./dataset/cropped/virat_kohli\\virat_kohli5.png',  
'./dataset/cropped/virat_kohli\\virat_kohli6.png',  
'./dataset/cropped/virat_kohli\\virat_kohli7.png',  
'./dataset/cropped/virat_kohli\\virat_kohli8.png',
```

```
class_dict = {}  
count = 0  
for celebrity_name in celebrity_file_names_dict.keys():  
    class_dict[celebrity_name] = count  
    count = count + 1  
class_dict
```

```
{'lionel messi': 0,  
 'maria_sharapova': 1,  
 'roger_federer': 2,  
 'serena_williams': 3,  
 'virat_kohli': 4}
```

Images in cropped folder can be used for model training. We will use these raw images along with wavelet transformed images to train our classifier. Let's prepare X and y now

```
X, y = [], []
for celebrity_name, training_files in celebrity_file_names_dict.items():
    for training_image in training_files:
        img = cv2.imread(training_image)
        scaled_raw_img = cv2.resize(img, (32, 32))
        img_har = w2d(img, 'db1', 5)
        scaled_img_har = cv2.resize(img_har, (32, 32))
        combined_img = np.vstack((scaled_raw_img.reshape(32*32*3,1), scaled_img_har.reshape(32*32*3,1)))
        X.append(combined_img)
        y.append(class_dict[celebrity_name])
```

```
len(X[0])
```

4096

```
32*32*3 + 32*32
```

4096

```
X[0]
```

```
array([[100],
       [129],
       [140],
       ...,
       [237],
       [234],
       [232]], dtype=uint8)
```

```
y[0]
```

0

```
X = np.array(X).reshape(len(X),4096).astype(float)
X.shape
```

(168, 4096)

Data cleaning process is done. Now we are ready to train our model

We will use SVM with rbf kernel tuned with heuristic finetuning

```
from sklearn.svm import SVC
from sklearn.preprocessing import StandardScaler
from sklearn.model_selection import train_test_split
from sklearn.pipeline import Pipeline
from sklearn.metrics import classification_report
```

```
X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=0)

pipe = Pipeline([('scaler', StandardScaler()), ('svc', SVC(kernel = 'rbf', C = 10))])
pipe.fit(X_train, y_train)
pipe.score(X_test, y_test)
```

0.8809523809523809

```
print(classification_report(y_test, pipe.predict(X_test)))
```

	precision	recall	f1-score	support
0	0.86	0.86	0.86	7
1	1.00	0.90	0.95	10
2	1.00	0.71	0.83	7
3	0.78	1.00	0.88	7
4	0.83	0.91	0.87	11
micro avg	0.88	0.88	0.88	42
macro avg	0.89	0.88	0.88	42
weighted avg	0.90	0.88	0.88	42

Let's use GridSearch to try out different models with different paramets. Goal is to come up with best modle with best fine tuned parameters

```
from sklearn import svm
from sklearn.ensemble import RandomForestClassifier
from sklearn.linear_model import LogisticRegression
from sklearn.pipeline import make_pipeline
from sklearn.model_selection import GridSearchCV
```

```
model_params = {
    'svm': {
        'model': svm.SVC(gamma='auto', probability=True),
        'params': {
            'svc__C': [1, 10, 100, 1000],
            'svc__kernel': ['rbf', 'linear']
        }
    },
    'random_forest': {
        'model': RandomForestClassifier(),
        'params': {
            'randomforestclassifier__n_estimators': [1, 5, 10]
        }
    },
    'logistic_regression': {
        'model': LogisticRegression(solver='liblinear', multi_class='auto'),
        'params': {
            'logisticregression__C': [1, 5, 10]
        }
    }
}
```

```

scores = []
best_estimators = {}
import pandas as pd
for algo, mp in model_params.items():
    pipe = make_pipeline(StandardScaler(), mp['model'])
    clf = GridSearchCV(pipe, mp['params'], cv=5, return_train_score=False)
    clf.fit(X_train, y_train)
    scores.append({
        'model': algo,
        'best_score': clf.best_score_,
        'best_params': clf.best_params_
    })
    best_estimators[algo] = clf.best_estimator_

df = pd.DataFrame(scores, columns=['model', 'best_score', 'best_params'])
df

```

C:\ProgramData\Anaconda3\lib\site-packages\sklearn\model_selection_search.py:84
 1: DeprecationWarning: The default of the `iid` parameter will change from True to False in version 0.22 and will be removed in 0.24. This will change numeric results when test-set sizes are unequal.

DeprecationWarning)

C:\ProgramData\Anaconda3\lib\site-packages\sklearn\model_selection_search.py:84
 1: DeprecationWarning: The default of the `iid` parameter will change from True to False in version 0.22 and will be removed in 0.24. This will change numeric results when test-set sizes are unequal.

DeprecationWarning)

C:\ProgramData\Anaconda3\lib\site-packages\sklearn\model_selection_search.py:84
 1: DeprecationWarning: The default of the `iid` parameter will change from True to False in version 0.22 and will be removed in 0.24. This will change numeric results when test-set sizes are unequal.

DeprecationWarning)

	model	best_score	best_params
0	svm	0.841270	{'svc__C': 1, 'svc__kernel': 'linear'}
1	random_forest	0.706349	{'randomforestclassifier__n_estimators': 10}
2	logistic_regression	0.809524	{'logisticregression__C': 1}

best_estimators

```

{'svm': Pipeline(memory=None,
  steps=[('standardscaler', StandardScaler(copy=True, with_mean=True, with_std=True)), ('svc', SVC(C=1, cache_size=200, class_weight=None, coef0=0.0, decision_function_shape='ovr', degree=3, gamma='auto', kernel='linear', max_iter=-1, probability=True, random_state=None, shrinking=True, tol=0.001, verbose=False))], 'random_forest': Pipeline(memory=None,
  steps=[('standardscaler', StandardScaler(copy=True, with_mean=True, with_std=True)), ('randomforestclassifier', RandomForestClassifier(bootstrap=True, class_weight=None, criterion='gini', max_depth=None, max_features='auto', max_leaf_nodes=None, min_impurity_decrease=0.0, min...obs=None, oob_score=False, random_state=None, verbose=0, warm_start=False))], 'logistic_regression': Pipeline(memory=None,
  steps=[('standardscaler', StandardScaler(copy=True, with_mean=True, with_std=True)), ('logisticregression', LogisticRegression(C=1, class_weight=None, dual=False, fit_intercept=True, intercept_scaling=1, max_iter=100, multi_class='auto', n_jobs=None, penalty='l2', random_state=None, solver='liblinear',

```

```
best_estimators['svm'].score(X_test,y_test)
```

```
0.9047619047619048
```

```
best_estimators['random_forest'].score(X_test,y_test)
```

```
0.5952380952380952
```

```
best_estimators['logistic_regression'].score(X_test,y_test)
```

```
0.9285714285714286
```

```
best_clf = best_estimators['svm']
```

```
from sklearn.metrics import confusion_matrix  
cm = confusion_matrix(y_test, best_clf.predict(X_test))  
cm
```

```
array([[ 7,  0,  0,  0,  0],  
       [ 0,  9,  0,  1,  0],  
       [ 0,  0,  5,  2,  0],  
       [ 0,  1,  0,  6,  0],  
       [ 0,  0,  0,  0, 11]], dtype=int64)
```

```
import seaborn as sn  
plt.figure(figsize = (10,7))  
sn.heatmap(cm, annot=True)  
plt.xlabel('Predicted')  
plt.ylabel('Truth')
```

```
Text(69.0, 0.5, 'Truth')
```

```
class_dict
```

```
{'lionel_messi': 0,  
 'maria_sharapova': 1,  
 'roger_federer': 2,  
 'serena_williams': 3,  
 'virat_kohli': 4}
```

Save the trained model

```
!pip install joblib  
import joblib  
# Save the model as a pickle in a file  
joblib.dump(best_clf, 'saved_model.pkl')
```

Requirement already satisfied: joblib in c:\programdata\anaconda3\lib\site-packages (0.15.1)
['saved_model.pkl']

Save class dictionary

```
import json  
with open("class_dictionary.json", "w") as f:  
    f.write(json.dumps(class_dict))
```