

# Comparing the performance of Emotion-Recognition Implementations in OpenCV, Cognitive Services, and Google Vision APIs

LUIS ANTONIO BELTRÁN PRIETO, ZUZANA KOMÍNKOVÁ OPLATKOVÁ

Department of Informatics and Artificial Intelligence

Tomas Bata University in Zlín

Nad Stráněmi 4511, 76005, Zlín

CZECH REPUBLIC

beltran\_prieto@fai.utb.cz

*Abstract:* - Emotions represent feelings about people in several situations. Various machine learning algorithms have been developed for emotion detection in a multimedia element, such as an image or a video. These techniques can be measured by comparing their accuracy with a given dataset in order to determine which algorithm can be selected among others. This paper deals with the comparison of three implementations of emotion recognition in faces, each implemented with specific technology. OpenCV is an open-source library of functions and packages mostly used for computer-vision analysis and applications. Cognitive services, as well as Google Cloud AI, are sets of APIs which provide machine learning and artificial intelligence algorithms to develop smart applications capable of integrate computer-vision, speech, knowledge, and language processing features. Three Android mobile applications were developed in order to test the performance between an OpenCV algorithm for emotion recognition, an implementation of Emotion cognitive service, and a Google Cloud Vision deployment for emotion-detection in faces. For this research, one thousand tests were carried out per experiment. Our findings show that the OpenCV implementation got the best performance, which can be improved by increasing the sample size per emotion during the training step.

*Key-Words:* - Emotion recognition, OpenCV, Fisherfaces, Cognitive Services, Cloud Vision, face detection

## 1 Introduction

Facial emotion recognition seeks to predict the real feeling that a person expresses based on facial images, with a wide range of possible applications, such as improving student engagement [1], building smart health environments [2], analyzing customers' feedback [3] and evaluating the quality of children's games [4], just to name a few.

Deep learning is a recent, revolutionary technique in machine learning which pursues the objective of bringing artificial intelligence to solve practical applications across different, diverse fields, such as recommender systems [5], plasma tomography reconstruction [6], facial age estimation [7], and neuroimaging [8], among others.

Recognizing faces within images and videos has been one of the challenges that deep learning has tested thoroughly, with significant performance and improvement. This progress has been achieved thanks to development of several techniques, such as Convolutional Neural Networks, Deep Belief Networks, and the availability of huge training datasets [9].

As a consequence, detecting the sentiment expressed by a person is the next step into facial

analysis. Recent research [10] has proven that emotion detection can be achieved by the usage of machine learning and artificial intelligence algorithms. While it is not an easy task, several open-source libraries and packages, such as OpenCV, TensorFlow, Theano, Caffe and CNTK (Microsoft Cognitive Toolkit) simplify the process of building deep-learning-based algorithms and applications. Emotions such as anger, disgust, happiness, surprise, and neutrality can be detected.

This paper is organized as follows. First, the problem formulation is mentioned along with a theoretical background on emotion recognition, OpenCV, Fisherfaces algorithm, Cognitive Services, Google Cloud AI, and the extended Cohn-Kanade (CK+) database is presented. Afterwards, the problem solution is described by explaining the methods and methodology that was used for this comparison in addition to the evaluation results. Finally, conclusions are discussed at the end of the paper.

## 2 Problem Formulation

The aim of this paper is to compare the performance of three emotion-recognition implementations. The

first application consists of Python code which uses OpenCV face-recognizer classes combined with Fisher Face technique. The second one is a C# application which sends requests to a Cognitive Services API. Finally, an Android application which makes Google's Cloud Vision API calls was used in the third application. 500 facial expressions from the Cohn-Kanade (CK+) dataset were examined by each program for evaluation purposes.

## 2.1 Background information

Emotions are strong feelings about people's situations and relationships with others. Most of the time, humans show how they feel by using facial expressions. Speech, gestures, and actions are also used to describe a person's current state.

Emotion recognition can be defined as the process of detecting the feeling expressed by humans from their facial expressions, such as anger, happiness, sadness, deceitfulness, and others. Even though a person can automatically identify facial emotions, machine learning algorithms have been developed for this purpose. Emotions play a key role in decision-making and human-behaviour, as many actions are determined by how a person feels at some point.

Typically, these algorithms use either a picture or a video (which can be considered as a set of images) as input, then they proceed to detect and focus their attention on a face and finally, specific points and regions of the face are analysed in order to detect the affective state.

Machine Learning (ML) algorithms, methods and techniques can be applied to detect emotions from a picture or video. For instance, a deep learning neural network can perform effective human activity recognition with the aid of smartphone sensors [11]. Moreover, a classification of facial expressions based on Support Vector Machines was developed for spontaneous behavior analysis.

## 2.2 OpenCV and Fisherfaces

OpenCV [12] is a free, yet powerful, open-source library developed by Intel Corporation which has been widely used in computer vision and machine learning tasks, such as image processing, real-time image recognition, and face detection. With more than 2500 optimized algorithms included, this library has been extensively used for research and commercial applications from both global and small entrepreneurs. OpenCV contains an optimized set of libraries written in C language, with bindings to

other languages and technologies, including Python, Android, iOS, and CUDA (for GPU fast processing), and wrappers in other languages, such as C#, Perl, Haskell, and others. Moreover, it works under Windows and Linux.

Among its capabilities, OpenCV contains a FaceRecognizer class which, as the name suggests, is helpful for face recognition tasks. There are three algorithms available for this purpose: Eigenfaces, Fisherfaces, and Local Binary Patterns Histograms. While the first technique considers a linear combination of facial features in order to maximize total variance in data, thus representing data in a powerful, but classless, way, Fisherfaces takes a Linear Discriminant Analysis approach in which class-specific dimensionality reduction is performed so the combination of features that separate the best classes is taken into account. If there exists any external source, such as light, which affects the representation of the image, the Eigenfaces technique is not able to accurately classify the faces. Fisherfaces, on its side, is not affected by this factor.

The algorithm goes as follows: First, let  $V = \{V_1, V_2, \dots, V_c\}$ ,  $V_i = \{v_1, v_2, \dots, v_N\}$  be a random vector with samples obtained from  $c$  classes. Then, the total mean,  $\mu$ , and the mean of class  $i$ ,  $\mu_i$ , where  $i \in \{1, 2, \dots, c\}$  are computed as described in equations (1) and (2). These values are used in equations (3) and (4) in order to calculate the Scatter matrices,  $S_B$  and  $S_W$ .

$$\mu = \frac{1}{N} \sum_{i=1}^N v_i \quad (1)$$

$$\mu_i = \frac{1}{|V_i|} \sum_{v_j \in V_i} v_j \quad (2)$$

$$S_B = \sum_{i=1}^c N_i (\mu_i - \mu) (\mu_i - \mu)^T \quad (3)$$

$$S_W = \sum_{i=1}^c \sum_{x_j \in x_i} (x_j - \mu_i) (x_j - \mu_i)^T \quad (4)$$

A projection  $W$  maximizes the class separability criterion by following equation (5):

$$W_{opt} = \arg \max_W \frac{|W^T S_B W|}{|W^T S_W W|} \quad (5)$$

The General Eigenvalue Problem solves this optimization task (6):

$$\begin{aligned} S_B v_i &= \lambda_i S_W v_i \\ S_W^{-1} S_B v_i &= \lambda_i v_i \end{aligned} \quad (6)$$

The rank of the scatter matrix,  $S_W$ , is at most ( $N$  samples –  $c$  classes). In problems such as pattern recognition tasks, the number of samples is smaller than the dimension of the data input. Thus,  $S_W$  becomes singular and can be solved using a linear discriminant analysis. At the end, the optimization problem is rewritten to equations (7) and (8):

$$W_{pca} = \arg \max_W |W^T S_T W| \quad (7)$$

$$W_{fld} = \arg \max_W \frac{|W^T W_{pca}^T S_B W_{pca} W|}{|W^T W_{pca}^T S_W W_{pca} W|} \quad (8)$$

And the transformation matrix  $W$ , projecting a sample into the  $(c - 1)$  dimensional space is given by equation (9):

$$W = W_{fld}^T W_{pca}^T \quad (9)$$

### 2.3 Cognitive Services

Cognitive Services [13] are a set of machine learning algorithms developed by Microsoft which are able to solve artificial intelligence problems in several fields, such as computer vision, speech recognition, natural language processing, machine learning search, and recommendation systems, among others. These algorithms can be consumed through Representational State Transfer (REST) calls over an Internet connection, allowing developers to use artificial intelligence research to solve problems. These services are open-source and can be consumed by many programming languages, including C#, PHP, Java, Python, and implemented in desktop, mobile, console, and web applications.

The Computer Vision API of Cognitive Services provides access to machine learning algorithms capable of performing image processing tasks. Either an image stream is uploaded or an image URL is specified to the service so the content can be analyzed for label tagging, image categorization, face detection, color extraction, text detection, and emotion recognition. Video is also supported as input. The Emotion API analyses the sentiment of a person in an image or video and returns the confidence level for eight emotions mostly understood by a facial expression, including anger, contempt, disgust, fear, happiness, neutrality, sadness and surprise.

### 2.4 Google Cloud AI

Google Cloud AI [14] is another collection of powerful machine learning services such as computer vision, speech recognition, text analysis, text translation and several others that can be used

to develop smart applications through REST requests by any programming language, in a similar way to the Cognitive Services API.

The Cloud Vision API of Google Cloud integrates several machine learning models capable of performing image classification tasks, such as image classification, object and face detection, optical character recognition (OCR), and even web detection, which searches for similar images on the Internet. An image, either from an URL or an array of bytes is provided to the service for the analysis to be performed. There are four sentiments detected by the API: joy, sorrow, anger, and surprise, along with five categorical label estimates representing the emotional confidence: very likely, likely, possible, unlikely, and very unlikely.

### 2.5 The Extended Cohn-Kanade database

The Cohn-Kanade AU-Coded Facial Expression Database [15][16] is a well-known repository of face images used for research purposes into the facial recognition field, with an increased interest in emotion detection research. An Extended version 2 of the database, also known as CK+, was developed in order to address some limitations of version 1, such as non-validation of emotion labels, the absence of a common performance metric against which to evaluate the latest algorithms and standard protocols for typical databases and quantitative meta-analysis. It contains the facial expressions of 210 adults between 18 and 50 years of age, from which 31% were male, 81% Euro-Americans, 13% Afro-Americans, and 6% from other groups. For each participant, 23 facial displays were performed. 593 sequences were labelled with a basic emotion from a pool of seven categories: anger, contempt, disgust, fear, happiness, sadness, and surprise.

## 3 Problem Solution

### 3.1 Methods and Methodology

The objective of this experiment is to compare the performance of three emotion-recognition implementations. The first analysis (Experiment A) is a Python-based application which makes use of the OpenCV machine learning algorithms for facial and emotion detection. The second study (Experiment B) is a C# mobile application which makes requests to a Cognitive Services API for emotion detection. Finally, the last test (Experiment C) consists of another C# mobile application which sends requests to a Google Cloud Vision API to

detect facial sentiments. In the three cases, the Extended Cohn-Kanade dataset of images was used as input for the analysis. All of the applications were developed by the authors for this experiment.

For Experiment A, we considered 327 sequences which actually show a relevant sequence of emotion, from a neutral feeling to the emotion itself. First step is then to obtain both the neutral and emotional images. From this subset, OpenCV library is used to detect the face on each picture by using a custom Haar-filter. Effective object detection is possible thanks to the Haar feature-based cascade classifiers proposed in [17], a machine learning based approach which takes advantage of cascade functions training from both positive and negative images. OpenCV already provides several Haar-filter functions; however, a cascade function of boosted classifiers was used for this experiment in order to detect the faces in the pictures. Thereafter, all images were standardized by converting them to grayscale and resized to the same dimensions. Table 1 presents how half of the subset, i.e., 327 images, were distributed among the different labelled emotions. The other half corresponds to 327 emotionless faces, i.e., showing neutrality. Then, we proceeded to randomly split the subset in two new sets: training set and prediction set. For training, we considered 523 images, which corresponds to 80% of the pictures. The remaining 20% was considered for the prediction set. For better evaluation purposes, 10 random training and classification sets were generated. The training process consists of getting the characteristics of each face along with the labelled emotion, i.e. the expression shown by the person. This data is used to create and train a Fisherface classifier. Then, evaluation of the classifier proceeds by comparing the outcome of its predict function of each face with the actual labelled emotion.

Experiment B starts with the 654 images extracted with the Python code from Experiment A. 10 random groups consisting each of 20% of the collection were generated in order to evaluate samples of the same size as in Experiment A. For each test, every face was submitted to the Cognitive Services API for its evaluation, as its training has already been developed by Microsoft. A C# mobile application was developed in Xamarin for running this experiment. The service returns a JSON content which contains the score for each emotion. Only the highest score, considered as the predicted facial expression detected by the service, was compared with the actual labelled emotion for evaluation purposes. Fig. 1 shows the application developed in C# for this experiment. It takes a previously

provided picture, then finds a face on it and finally submits a request in order to detect the emotions expressed by the person in the photo. The analysis is performed by the Emotion Cognitive Service.

**Table 1.** Sample size of each labelled emotion

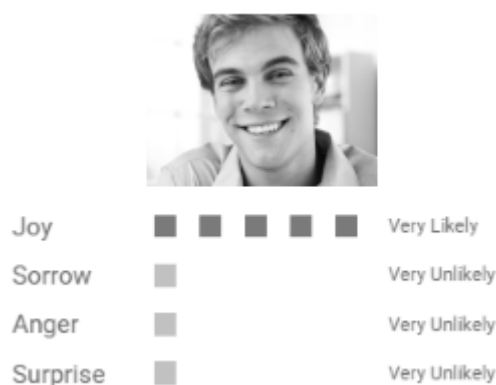
Labelled Emotion	Number of faces
Anger	45
Contempt	18
Disgust	58
Fear	25
Happiness	69
Sadness	28
Surprise	84



**Fig. 1.** Analysis of emotions detected on a picture by the Emotion Cognitive Service.

Only images from the following labelled subsets were considered for Experiment C: Surprise, Anger, Happiness, and Sadness, due to the fact that the Cloud Vision API identifies four sentiments: Surprise, Anger, Joy, and Sorrow, respectively. A total of 226 faces were considered as a result for this experiment. 10 random groups consisting each of 20% of the collection were generated for sentiment detection. For every test, each face was included in a request to the Cloud Vision API for its analysis, as the model training has already been developed by Google. A C# mobile application mobile application was developed with Xamarin technology for running this experiment. In a similar way to its Microsoft counterpart, the service returns a JSON

string, which contains the dominant emotion along with its likelihood estimation. This observed expression was compared with the actual labelled emotion for evaluation purposes. Fig. 2 shows the software application developed in C# for this experiment. It takes a previously provided picture, then finds a face on it and finally submits a request in order to detect the emotions expressed by the person in the photo. The analysis is performed by the Cloud Vision service.



**Fig. 2.** Analysis of emotions detected on a picture by the Cloud Vision Service.

### 3.2 Results and Discussion

After running each of the 10 tests from Experiment A, the results which are presented in Table 2 were obtained. An average of 75.87% correct predictions was calculated as an outcome. Likewise, Table 3 illustrates the outcome of each test in Experiment B. As a result, a 68.93% average efficiency was accomplished after running 10 tests of this implementation. Lastly, Table 4 displays the observations for Experiment C after a 10-test evaluation. A 72.74% accuracy average was achieved subsequently.

The findings of the experiments show that the Python-based implementation with OpenCV using Fisherfaces proved to be more accurate than both the Cognitive Services and the Cloud Vision API implementations in C# by approximately a 7% difference. For experiment A, several of the mistakes occurred when trying to predict emotions with low occurrences in the dataset, such as fear and contempt mistakenly classified as neutral expressions. Moreover, there were a few errors when trying to predict a neutral face, most of the time identified as a sad face. Regarding experiment B, neutral images were wrongly identified as either contempt or sadness emotions; however, by looking closely to the scores obtained by the Cognitive Services, a minimum difference between the wrong

prediction and the actual emotion was detected. Thus, in most cases, the second-best prediction was correct. However, for evaluation purposes, this was considered as an error. For Experiment C, several faces were detected by the Cloud Vision API as sad expressions (with a low likelihood, however enough to be considered as the dominant emotion) while the real expressed sentiment was anger.

**Table 2.** Evaluation results of experiment A

Test Number	Correct (%)	Incorrect (%)
1	103 (78.62%)	28 (21.37%)
2	100 (76.33%)	31 (23.66%)
3	98 (74.80%)	33 (25.19%)
4	104 (79.38%)	27 (20.61%)
5	96 (73.28%)	35 (26.71%)
6	95 (75.25%)	36 (27.48%)
7	99 (75.57%)	32 (24.42%)
8	100 (76.33%)	31 (23.66%)
9	97 (74.04%)	34 (25.95%)
10	102 (77.86%)	29 (22.13%)

**Table 3.** Evaluation results of experiment B

Test Number	Correct (%)	Incorrect (%)
1	89 (67.93%)	42 (32.06%)
2	95 (72.51%)	36 (27.48%)
3	86 (65.64%)	45 (34.35%)
4	91 (69.46%)	40 (30.53%)
5	88 (67.17%)	43 (32.82%)
6	90 (68.70%)	41 (31.29%)
7	95 (72.51%)	36 (27.48%)
8	93 (70.99%)	38 (29.00%)
9	84 (64.12%)	47 (35.87%)
10	92 (70.22%)	39 (29.77%)

**Table 4.** Evaluation results of experiment C

Test Number	Correct (%)	Incorrect (%)
1	95 (72.51%)	36 (27.48%)
2	90 (68.70%)	41 (31.29%)
3	103 (78.62%)	28 (21.37%)
4	91 (69.46%)	40 (30.53%)
5	97 (74.04%)	34 (25.95%)
6	92 (70.22%)	39 (29.77%)
7	94 (71.75%)	37 (28.24%)
8	96 (73.28%)	35 (26.71%)
9	99 (75.57%)	32 (24.42%)
10	96 (73.28%)	35 (26.71%)

#### 4 Conclusion

The objective of this experiment was to compare the performance of three different implementations of emotion-recognition applications by using OpenCV and Python with a Fisherface technique in the first case, while considering a C#-based solution which makes requests to a Cognitive Services API for Emotion detection for the second solution. A third analysis included a C#-based software which sends API calls to a Cloud Vision service for sentimental analysis of pictures. While the first implementation got the best results, the performance could be improved either by increasing the sample size of those emotions with few faces, so the training phase gets benefited, or by removing them from the subset, as not enough cases were collected.

#### Acknowledgement

This work was supported by the Ministry of Education, Youth and Sports of the Czech Republic within the National Sustainability Programme project No. LO1303 (MSMT-7778/2014) and also by the European Regional Development Fund under the project CEBIA-Tech No. CZ.1.05/2.1.00/03.0089, further it was supported by Grant Agency of the Czech Republic—GACR 588 P103/15/06700S and by Internal Grant Agency of Tomas Bata University in Zlin under the project No. IGA/CebiaTech/2017/004. L.A.B.P author also thanks the doctoral scholarship provided by the National Council for Science and Technology (CONACYT) and the Council for Science and Technology of the State of Guanajuato (CONCYTEG) in Mexico.

#### References:

- [1] Garn AC, Simonton K, Dasingert T, Simonton A. Predicting changes in student engagement in university physical education: Application of control-value theory of achievement emotions, *Psychology of Sport and Exercise*, Vol.29, pp. 93-102.
- [2] Fernandez-Caballero A, Martinez-Rodrigo A, Pastor JM, Castillo JC, Lozano-Monator E, Lopez MT, Zangroniz R, Latorre JM, Fernandez-Sotos A, Smart environment architecture for emotion detection and regulation, *Journal of Biomedical Informatics*, Vol.64 pp-57-73.
- [3] Felbermayr A, Nanopoulos A, The Role of Emotions for the Perceived Usefulness in Online Customer Reviews, *Journal of Interactive Marketing*, Vol.36, pp. 60-76.
- [4] Gennari R, Melonio A, Raccanello D, Brondino M, Doderio G, Pasini M, Torello S, Children's emotions and quality of products in participatory game design, *International Journal of Human-Computer Studies*, Vol.101, pp. 45-61.
- [5] Wei J, He J, Chen K, Zhou Y, Tang Z, Collaborative filtering and deep learning based recommendation system for cold start items, *Expert Systems with Applications*, Vol.69, pp. 29-39.
- [6] Matos FA, Ferreira DR, Carvalho PJ, Deep learning for plasma tomography using the bolometer system at JET, *Fusion Engineering and Design*, Vol.114, pp. 18-25.
- [7] Liu H, Lu J, Feng J, Zhou J, Group-aware deep feature learning for facial age estimation, *Pattern Recognition*, Vol.66, pp. 82-94.
- [8] Vieira S, Pinaya WHL, Mechelli A, Using deep learning to investigate the neuroimaging correlates of psychiatric and neurological disorders: Methods and applications, *Neuroscience and Biobehavioral Reviews*, Vol.74, pp. 58-75.
- [9] Parkhi OM, Vedaldi A, Zisserman A, Speeding up Convolutional Neural Networks with Low Rank Expansions, *Proceedings of the British Machine Vision Conference*, BMVA Press, 2014.
- [10] Yogesh CK, Hariharan M, Ngadiran R, Adom AH, Yaacob S, Polat K, Hybrid BBO\_PSO and higher order spectral features for emotion and stress recognition from natural speech, *Applied Soft Computing*, Vol.56, pp. 217-232.
- [11] Ronao CA, Cho S, Human activity recognition with smartphone sensors using deep learning

- neural networks, *Expert Systems with Applications*, Vol.59, pp. 235-244.
- [12] OpenCV library. <http://opencv.org> [Online: accessed 01-Jun-2017]
- [13] Cognitive Services – Intelligence Applications. <http://microsoft.com/cognitive> [Online: accessed 03-Jun-2017]
- [14] Google Cloud Machine Learning at Scale <https://cloud.google.com/products/machine-learning/> [Online: accessed 21-Ago-2017]
- [15] Kanade T, Cohn JF, Tian Y, Comprehensive database for facial expression analysis., *Proceedings of the Fourth IEEE International Conference on Automatic Face and Gesture Recognition (FG'00)*, pp. 46-53.
- [16] Lucey P, Cohn JF, Kanade T, Saragih J, Ambadar Z, Matthews I, The Extended Cohn-Kanade Dataset (CK+): A complete expression dataset for action unit and emotion-specified expression, *Proceedings of the Third International Workshop on CVPR for Human Communicative Behavior Analysis (CVPR4HB 2010)*, pp. 94-101.
- [17] Turk M, Pentland A, Eigenfaces for Recognition, *Journal of Cognitive Neuroscience*, Vol.3, No. 1, pp. 71-86.