

Cloud-based Facial Emotion Recognition for Real-Time Emotional Atmosphere Assessment during a Lecture

Peter TAKÁČ*, Marián MACH**, Peter SINČÁK***

Dept. of Cybernetics and Artificial Intelligence, FEI TU of Košice, Slovak Republic

* peter.takac.3@tuke.sk, ** marian.mach@tuke.sk, *** peter.sincak@tuke.sk

Abstract—These instructions give you basic guidelines for preparing camera-ready papers for conference proceedings.

I. INTRODUCTION

Emotions in inter-human communication are being proven every day to be one of the main drivers of the communication topic change, response formulation and decision making. However, the emotion recognition in face-to-face communication is relatively easy task for the humans for the machine it is a very difficult task while the modality of the emotion recognition depends on multiple factors during the communication [1]. In contrast with humans who can perceive the communication modalities with ease the machine has to split the task of multi-modal emotion recognition into subtasks which should be easier to carry out. Therefore, we divide emotion recognition to groups which are defined by the source of the expressed emotion. One of the most frequently processed groups of emotions are facial emotions while they offer a kind of universality across multiple cultures. This universality theory was studied since the early 1970s when Ekman and Friesen [2] performed extensive studies of human facial expressions, providing evidence to support the theory. By relying on this theory we are able to create systems that are able to recognize facial emotions within the given range of intensity for each human.

The precision of the facial recognition is prone to change considering the amount of facial features which are either not present, obstructed, not tracked, or not visible. Hence, a facial emotion recognition system is highly dependent on the quality of acquired input data i.e. the human face. For a high quality input we can consider a picture of a face that is acquired in high resolution with zero or close to zero compression, the face is still, upright frontal with no obstructions (glasses, patches, not covered by any object) and is not influenced by lighting conditions i.e. absence of underexposed or overexposed areas [3]. However, the acquisition of such inputs is usually difficult and the images often suffer from deficiencies most commonly related to lighting, motion, or poor resolution of the sensor. Therefore, a preprocessing step is needed which according to computer vision can be computationally demanding, and if the desired number of faces in the processed picture raises also the computational demands raise exponentially. These difficulties can be usually solved by use of hardware with better stats by the means of computational performance where a cloud-based solutions seem to be one of the best performing due to scalability of cloud's resources.

Facial emotion recognition in the role of social interaction between machines and humans can be considered as a tool for additional information acquisition in feedback driven interaction where the machine is acting according to human's reaction or behavior and the emphasis is on information transfer [4]. This additional information is important for emotional atmosphere evaluation which leads to better adaptation to the interaction topic instead of forcibly changing the topic.

The use of a machine-based emotion recognition is advantageous when we switch from emotion recognition in a single face to multiple faces. The difference between machines and humans in this case grows even further if the environment is rapidly changing, for example during a group discussion. This advantage is based on the ability of the machine to work in parallel and therefore the machine can recognize emotions of multiple humans at once (in case every face is located on the same image or some mechanism is able to join the scanned faces). This advantage can be used by a human user who needs to be able to track multiple people at once. In this paper we propose a task during which such system can be utilized i.e. emotional atmosphere assessment during a lecture. This scenario is common in almost any field where a presenter or a lecturer wants to achieve the best results in responsiveness of his audience during his speech or presentation.

II. EMOTION RECOGNITION IN CLOUD ENVIRONMENT

For facial emotion recognition there is a need for an exact universal model, which ensures correct emotion classification. The means by which one emotion is distinguished from another have been researched from two fundamental viewpoints. That emotions are discrete and fundamentally different constructs [5], or that emotions can be characterized on a dimensional basis [6].

A. Discrete emotional categories

The discrete emotion theory [7], claims that there is a small number of core emotions which are biologically determined emotional responses whose expression and recognition is fundamentally the same for all individuals regardless of ethnic or cultural differences. One of the most popular example of discrete emotion categories is Ekman's model [8] in which they state that there are six basic emotions (anger, disgust, fear, happiness, sadness, and surprise). Furthermore, it is explained that there are particular characteristics attached to each of these

emotions, allowing them to be expressed in various degree. However, each recognized facial emotion acts as a discrete category rather than an individual emotional state which corresponds to the problem of the multimodal character of human emotions. Nevertheless, for purposes of approximate reaction emotion assessment the facial emotions are suitable and while we were interested in an immediate feedback during a lecture we used a discrete emotional model to recognize emotions of students. For these purposes we use the Microsoft Emotion API [9] which is available as cloud service on the Microsoft AZURE Cloud.

The Microsoft Emotion API uses a modified version of Ekman's emotion model to which they added two emotions namely contempt and neutral. Furthermore, we have created a simple verification procedure for determination of the service's performance using contingency tables and basic classification performance measures

$$Precision = \frac{TP}{TP + FP} \quad Recall = \frac{TP}{TP + FN}$$

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$

$$Error\ rate = \frac{FP + FN}{TP + TN + FP + FN}$$

Where TP is true positive, TN is true negative, FP is false positive and FN is false negative.

$$Random\ accuracy = \frac{\sum_{i=1}^M (\sum_{j=1}^M x_{ij} \cdot \sum_{k=1}^M x_{ki})}{(\sum_{i=1}^M (\sum_{j=1}^M x_{ij}))^2}$$

Where M represents the number of classes.

$$Cohen's\ kappa = \frac{Accuracy - Random\ accuracy}{1 - Random\ accuracy}$$

B. Dimensional models of emotions

The dimensional models of emotions attempt to conceptualize human emotions by defining their position in a multidimensional model. These models suggest that a common and interconnected neurophysiological system is responsible for all affective states [10]. In [10] the authors also present a circumplex model of emotions which suggests that emotions are distributed in a two-dimensional circular space, containing arousal and valence dimensions (see Figure 1.). Besides facial emotion recognition the circumplex models are broadly used for words emotion stimuli testing and affective state labeling [11].

For the emotional atmosphere evaluation we propose to use a modified version of Russell's circumplex model which serves for overall aggregated emotion specification where the modification lies in assumption that we consider arousal to be a measure of emotion intensity, and is used for further specification of acquired emotions according to

previous emotion measurements. For the conversion to this model we use the following formulas.

$$EEM_t = \sum_{i=1}^N \left(\frac{\sum_{k=1}^K p_{ik}}{K} \right) - \sum_{j=1}^M \left(\frac{\sum_{l=1}^K n_{jl}}{K} \right).$$

Where EEM_t (Ekman's Emotional Model) refers to aggregated emotions represented by Ekman's model; N is the number of positive emotions, and M is the number of negative emotions in the used model; K is the number of overall recognized faces from the current picture; p_{ik} refers to the value of positive emotion, and n_{jl} refers to the value of negative emotion.

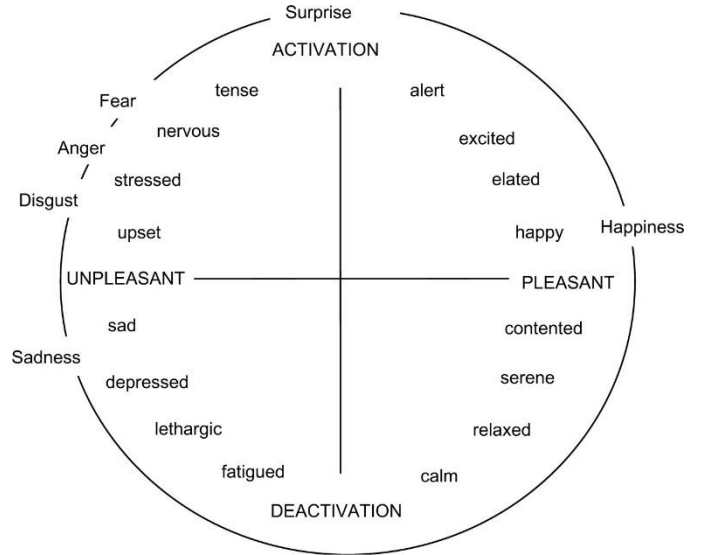


Figure 1. An example of representation of Russell's circumplex model of emotions with the marked position of Ekman's emotions. The horizontal axis represents the valence values and the vertical axis represents the arousal values.

Using the Microsoft Emotion API we can recognize only two emotions which we consider being positive (happiness and surprise), five emotions considered being negative (anger, contempt, disgust, fear, sadness), and the neutral emotion expression which is excluded from the evaluation process while it's contribution to aggregated emotions is zero and is only used for aggregated emotion intensity.

Next, we use a classic feature scaling normalization method for input normalization to $(-1,1)$ interval

$$Norm(x) = a + \frac{(x - x_{max})(b - a)}{x_{max} - x_{min}}$$

Where a is the lowest value of normalization interval ($a = -1$) and b is the highest value ($b = 1$); x_{min} and x_{max} are the minimum and maximum numbers of faces available to recognize. The normalization in our proposed solution is crucial while we wanted to be able to get an assessment of emotional atmosphere for various numbers of faces.

	A	C	D	F	H	N	SA	SU	SUM	Recall
A	0	1	0	0	0	23	6	0	30	0
C	0	0	0	0	0	0	0	0	0	0
D	0	0	2	0	0	10	17	0	29	0.07
F	0	0	0	6	0	13	11	2	32	0.19
H	0	0	0	0	28	3	0	0	31	0.9
N	0	0	0	0	0	30	0	0	30	1
SA	0	0	0	0	1	12	16	0	29	0.55
SU	0	0	0	0	2	5	0	22	29	0.76
SUM	0	1	2	6	31	96	50	24	210	
Precision	0	0	1	1	0.9	0.3	0.3	0.92		

A - anger, C - contempt, D - disgust, F - fear, H - happiness,
N - neutral, SA - sadness, SU - surprise

The results of the performance measures are following: *accuracy* = **49,5%**, *error rate* = **50,5%**, *random accuracy* = **14,14%**, *Cohen's kappa* = **41,2%**, and the *response precision* = **98,6%**. Despite the results show only nearly 50% accuracy the prediction of four out of seven emotions report 90-100% precision. This result also influences the Cohen's kappa which shows significantly higher value than the previous dataset results. We assume that the model which is used for emotion recognition in Emotion API was not thoroughly trained on databases such as JAFFE, but we would rate the ability to recognize emotions in cross-cultural faces as sufficient.

The last tested database was the Karolinska Directed Emotional Faces database. We have divided the testing of this database in three parts while we wanted to measure the accuracy of recognition on rotated faces. The KDEF database consists of 4900 faces where each face is labeled with 7 out of 8 emotions of the emotion vector of the Emotion API.

TABLE 3.
CONTINGENCY TABLE OF THE RESULTS FROM THE
COMPARISON WITH THE HALF-LEFT FACES FROM THE KDEF
DATABASE.

	A	C	D	F	H	N	SA	SU	SUM	Recall
A	6	0	5	0	1	30	1	0	43	0.14
C	0	0	0	0	0	0	0	0	0	0
D	5	0	52	0	6	5	3	0	71	0.73
F	1	0	4	2	5	14	9	9	44	0.05
H	1	0	0	0	61	0	0	0	62	0.98
N	0	0	0	0	0	47	0	0	47	1
SA	0	0	2	0	1	13	31	0	47	0.66
SU	0	0	0	0	0	2	0	4	6	0.67
SUM	13	0	63	2	74	111	44	13	320	
Precision	0.46	0	0.83	1	0.82	0.42	0.71	0.31		
A - anger, C - contempt, D - disgust, F - fear, H - happiness, N - neutral, SA - sadness, SU - surprise										

TABLE 4.
CONTINGENCY TABLE OF THE RESULTS FROM THE
COMPARISON WITH THE HALF-RIGHT FACES FROM THE
KDEF DATABASE.

	A	C	D	F	H	N	SA	SU	SUM	Recall
A	7	0	6	0	0	29	3	0	45	0.16
C	0	0	0	0	0	0	0	0	0	0
D	3	0	50	0	1	6	13	0	73	0.69
F	1	0	4	2	4	5	13	10	39	0.05
H	0	0	0	0	52	0	0	0	52	1
N	0	0	0	0	1	45	0	0	46	0.98
SA	1	0	0	0	1	9	36	0	47	0.77
SU	0	0	0	0	0	7	0	7	14	0.5
SUM	12	0	60	2	59	101	65	17	316	
Precision	0.58	0	0.83	1	0.88	0.45	0.55	0.41		
A - anger, C - contempt, D - disgust, F - fear, H - happiness, N - neutral, SA - sadness, SU - surprise										

TABLE 5.
CONTINGENCY TABLE OF THE RESULTS FROM THE
COMPARISON WITH THE UPRIGHT FRONTAL FACES FROM
THE KDEF DATABASE.

	A	C	D	F	H	N	SA	SU	SUM	Recall
A	63	1	9	1	0	55	10	0	139	0.45
C	0	0	0	0	0	0	0	0	0	0
D	9	0	101	1	1	3	24	0	139	0.73
F	0	3	7	26	3	17	32	47	135	0.19
H	0	0	0	0	136	0	0	0	136	1
N	0	0	0	0	0	134	2	0	136	0.99
SA	0	0	0	0	1	19	117	0	137	0.85
SU	0	0	0	0	3	10	0	126	139	0.91
SUM	72	4	117	28	144	238	185	173	961	
Precision	0.88	0	0.86	0.93	0.94	0.56	0.63	0.73		
A - anger, C - contempt, D - disgust, F - fear, H - happiness, N - neutral, SA - sadness, SU - surprise										

The real and computed values of the emotion assessment can be inspected in tables TABLE 3., TABLE 4., and TABLE 5. First results we have evaluated on this dataset were the results over rotated faces. There are only two types of rotations with the angle of 45 degrees either clockwise (half-right) or counterclockwise (half-left) from the straight upright position. The resultant measurements for the rotated faces were the following, half-left: *accuracy* = 63,4%, *error rate* = 36,6%, *random accuracy* = 16,7%, *Cohen's kappa* = 55,6%, *response precision* = 40%; half-right: *accuracy* = 63%, *error rate* = 37%, *random accuracy* = 16%, *Cohen's kappa* = 55,9%, *response precision* = 39,5%. As we can see the Emotion API performs with over 80% precision in recognition of three emotions (disgust, fear, and happiness) which show the highest difference from the other emotions. However, response precision is low due to poor face detection rate and therefore we would propose to use the Emotion API only for straight upright faces.

The results of straight upright faces recognition are following: *accuracy* = 73,2%, *error rate* = 26,8%, *random accuracy* = 14,23%, *Cohen's kappa* = 68,7%, *response precision* = 98%. This subset of images showed the highest accuracy throughout the tests. As we can see in the contingency table almost every emotion has precision over 70% which is sufficient for an approximate facial emotion assessment. We consider the results of this test as a success while the accuracy of this solution depends on specifying one exact emotion from the emotion vector which on the other hand returns approximate assessment of appearance of multiple emotions and the accuracy can be increased by definition of co-occurrence of specific emotions.

B. Emotional atmosphere assessment during a lecture

The next experiment we carried out was the emotional atmosphere assessment during a lecture. For this test, we created a setup with a Kinect V2 sensor which served as an input device for image acquisition and was connected to a laptop which ensured connection to the Microsoft AZURE Cloud over the internet. The internet connection we used was a classic 4G LTE with the maximum transfer

rate of 150 Mbit/s. The frequency of emotional measurement was every 3000 milliseconds while the pictures taken were in FullHD resolution (1920x1080 pixels) and this was the highest possible frequency allowed to upload these pictures in sequence using the 4G network.

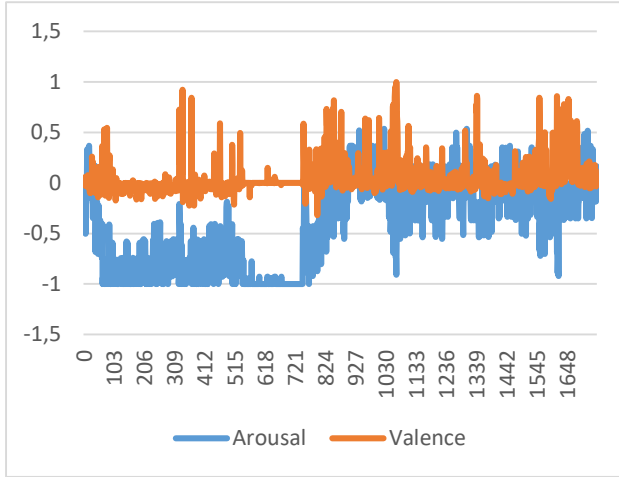


Figure 2. Graph of emotional assessment during a lecture, where the horizontal axis holds the time factor measured in frames acquired every 3000 milliseconds.

The students attending the lecture were monitored throughout the lecture, but for this scenario the lecturer had no information about the emotional assessment of the atmosphere. This test served as verification of the usability of Emotion API where the number of faces in image/frame is higher than one. The results of the experiment are shown in Figure 2. These results show various emotional assessments of students performing various types of actions. We divided the lecture into four parts: introduction, test, and the lecture.

During the introduction part, the students were informed about the fact that they will be monitored and were asked if they agree with it. They were also informed, about the upcoming test. This part was short around 80 frames i.e. 4 minutes. As we can see there is a low starting arousal value where the students were acquiring their seats and were not looking towards the camera, therefore, almost no emotion could be detected.

In the next part of the lecture, the students were writing a test which can be also clearly seen from the graph. The arousal values are low which indicate that the students were fully focusing on the test. We put a small joke at the beginning of each test to verify if the valence values will rise. This was confirmed by the risen valence values around 80 - 100 frames. We have given the students 25 minutes for the test (from 100 to 600 frames) and then they were instructed to leave the room for a 10-minute pause (600 to 800 frames).

After the test they had an ordinary lecture where we can see multiple changes in the arousal and valence values. This part took around 1100 frames which is around 55 minutes. During the whole experiment, we have evaluated a total of 1752 frames, transferred around 7 GB of data, recognized and evaluated emotions from around 21 000

faces. During the test the overall arousal value was $-0,85$ which indicates that most of the students were fully focused on the test, and the overall valence value was $-0,01$ which meant a neutral attitude towards writing the test. Outside the test the overall arousal value was $-0,11$ picturing slightly less attention paid to the lecture, and the overall valence of $0,07$ meant that the students had almost the same attitude as during the test.

IV. CONCLUSION

In this paper, we have reviewed the significance of machines in helping the human to recognize multiple emotions at once what can help the human user to acquire additional information about the environment which is hardly accessible or completely inaccessible. Next, we have also verified the performance of the Microsoft Emotion API where the evaluated results show sufficient performance according to tested databases. However, we conclude that the experiment does not have to necessarily picture the whole potential of this cloud-based emotion recognition while it is not certain if a single emotion label for a face is appropriate and we assume that in some extent also other emotions can play a role in the overall emotion assessment.

Lastly, we have created a solution for acquiring an emotional atmosphere assessment during a lecture using the mentioned Microsoft Emotion API. For this experiment, we propose a conversion method for converting the emotional vector of the Ekman's emotional model to a more feasible model for emotional atmosphere assessment – the Russell's circumplex model. We have prepared an experiment during which we wanted to prove the efficiency of conversion in combination with the Microsoft Emotion API. We have tested the solution during a lecture at our university where the students who were in the positions of the audience were instructed to perform various tasks. The experiment's results proved high correlation with the student's activities and therefore we claim that the proposed solution is usable for an emotional atmosphere assessment task.

Improvements in the verification can be made by acquiring a face database which would contain multiple values or occurrences of emotions. With such database, the verification procedure can be enhanced and from our view the resultant comparison could give higher accuracies. Furthermore, the conversion method which is used to convert the emotional vector of the Ekman's model to a Russell's circumplex model can be further enhanced by taking in count the degree of uncertainty of emotions. To work with the uncertainty we can use the representation of emotions or aggregated emotions in fuzzy logic. Therefore, we see the possibility of creation and application of fuzzy atmosphere assessment for various task concerning lecturing, inter-human, and human-machine communication i.e. social robotics.

ACKNOWLEDGEMENT

Research supported by the National Research and Development Project Grant 1/0773/16 2016 – 2019

“Cloud Based Artificial Intelligence for Intelligent Robotics” and by the Slovak Research and Development Agency under the contract No. APVV-015- 0731 and research project is supported from 07-2016 to 06-2019..

REFERENCES

- [1] A. James, N. Sebe, “Multimodal human-computer interaction: A survey.” *Computer vision and image understanding* 108, no.1, pp 116-134, 2007.
- [2] P. Ekman, “Darwin and facial expression: A century of research in review.” Ishk, 2006.
- [3] S. Agrawal, P. Khatri, “Facial expression detection techniques: based on Viola and Jones algorithm and principal component analysis.” *2015 Fifth International Conference on Advanced Computing & Communication Technologies (ACCT)*, pp. 108-112, IEEE, 2015.
- [4] R. W. Picard, “Affective Computing.” *Vol. 252*. MIT press, Cambridge, 1997.
- [5] G. Colombetti, “From affect programs to dynamical discrete emotions.” *Philosophical Psychology*, pp. 407-425, 2009.
- [6] W. M. Wundt, “Grundzüge der physiologischen Psychologie.” 1874.
- [7] C. E. Izard, “Emotion theory and research: Highlights, unanswered questions, and emerging issues.” *Annual review of psychology*. 2009.
- [8] P. Ekman, V. W. Friesen, “Constants across cultures in the face and emotion.” *Journal of personality and social psychology*, 1971.
- [9] <https://www.microsoft.com/cognitive-services/en-us/emotion-api>.
- [10] J. Posner, J. Russell, B.S. Peterson, “The cognitive development, and psychopathology.” *Development and psychopathology*, pp. 715-734, 2005.
- [11] N. A. Remington, L. R. Fabrigar, P. S. Visser, “Reexamining the circumplex model of affect.” *Journal of personality and social psychology*, 2000.
- [12] <https://azure.microsoft.com>.
- [13] <http://www.asp.net/web-api>.
- [14] M. M. Nordstrom, M. Larsen, J. Sierakowski, M. B. Stegmann, “The IMM Face Database – An Annotated Dataset of 240 Face Images.” *Informatics and Mathematical Modeling*, Technical University of Denmark, 2004.
- [15] M. J. Lyons, S. Akemastu, M. Kamachi, J. Gyoba, “Coding Facial Expressions with Gabor Wavelets.” *3rd IEEE International Conference on Automatic Face and Gesture Recognition*, pp. 200-205, 1998.
- [16] D. Lundqvist, A. Flykt, A. Öhman, “The Karolinska Directed Emotional Faces – KDEF.” *CD ROM from Department of Clinical Neuroscience*, Psychology section, Karolinska Institute, ISBN 91-630-7164-9, 1998.