

AVEC 2013 – The Continuous Audio/Visual Emotion and Depression Recognition Challenge

Michel Valstar
University of Nottingham
Mixed Reality Lab

Florian Eyben
TU München
MISP Group, MMK

Sebastian Schnieder
University of Wuppertal
Schumpeter School of
Business and Economics

Björn Schuller^{*}
TU München
MISP Group, MMK

Bihan Jiang
Imperial College London
Intelligent Behaviour
Understanding Group

Roddy Cowie
Queen's University
School of Psychology

Kirsty Smith
University of Nottingham
Mixed Reality Lab

Sanjay Bilakhia
Imperial College London
Intelligent Behaviour
Understanding Group

Maja Pantic[†]
Imperial College London
Intelligent Behaviour
Understanding Group

ABSTRACT

Mood disorders are inherently related to emotion. In particular, the behaviour of people suffering from mood disorders such as unipolar depression shows a strong temporal correlation with the affective dimensions valence and arousal. In addition, psychologists and psychiatrists take the observation of expressive facial and vocal cues into account while evaluating a patient's condition. Depression could result in expressive behaviour such as dampened facial expressions, avoiding eye contact, and using short sentences with flat intonation. It is in this context that we present the third Audio-Visual Emotion recognition Challenge (AVEC 2013). The challenge has two goals logically organised as sub-challenges: the first is to predict the continuous values of the affective dimensions valence and arousal at each moment in time. The second sub-challenge is to predict the value of a single depression indicator for each recording in the dataset. This paper presents the challenge guidelines, the common data used, and the performance of the baseline system on the two tasks.

^{*}The author is further affiliated with Imperial College London, Department of Computing, London, U.K.

[†]The author is further affiliated with Twente University, EEMCS, Twente, The Netherlands.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

AVEC'13, October 21, 2013, Barcelona, Spain.

Copyright © 2013 ACM 978-1-4503-2395-6/13/10...\$15.00.

<http://dx.doi.org/10.1145/2512530.2512533>.

Categories and Subject Descriptors

J [Computer Applications]: Miscellaneous; D.2.8 [Software Engineering]: Metrics—*complexity measures, performance measures*

Keywords

Affective Computing, Emotion Recognition, Speech, Facial Expression, Challenge

1. INTRODUCTION

According to European Union Green Papers dating from 2005 [15] and 2008 [16], mental health problems affect one in four citizens at some point during their lives. As opposed to many other illnesses, mental ill health often affects people of working age, causing significant losses and burdens to the economic system, as well as the social, educational, and justice systems. It is therefore somewhat surprising that despite the scientific and technological revolutions of the last half century remarkably little innovation has occurred in the clinical care of mental health disorders.

Affective Computing and Social Signal Processing are two of the more recent technological revolutions that promise to change this situation. Affective Computing is the science of automatically analysing affect and expressive behaviour [20]. By their very definition, mood disorders are directly related to affective state. Social Signal Processing addresses all verbal and non-verbal communicative signalling during social interactions, be they of an affective nature or not [27]. Although the measurement and assessment of behaviour is a central component of mental health practice it is severely constrained by individual subjective observation and lack of any real-time naturalistic measurements. It is thus only logical that researchers in affective computing and social signal processing, which aim to quantify aspects of expressive behaviour such as facial muscle activations and speech rate, have started looking at ways in which their communities can help mental health practitioners.

In the case of depression, which is the focus of AVEC 2013, the clinician-administered Hamilton Rating Scale for Depression [13] is the current gold standard to assess severity [2, 30], whereas the gold-standard for diagnosis is the Structured Clinical Interview for DSM-IV (SCID) [9]. The Hamilton scale is not free to use, but other self report measures are. The frequently-used Beck Depression Inventory [3] is one of them, and is the one used to obtain the ground truth measure for AVEC. All of these instruments pay little or no attention to observational behaviour. In part for that reason, social signal processing and affective computing could make significant contribution by achieving an objective, reliable method to incorporate behaviour into clinical assessment.

In the first published efforts towards this, the University of Pennsylvania has already applied a basic facial expression analysis algorithm to distinguish between patients with Schizophrenia and healthy controls [28, 14]. Besides diagnosis, affective computing and social signal processing would also allow quantitative monitoring of the progress and effectiveness of treatment. Early studies that addressed the topic of depression are e.g. [28, 4].

More recently, Girard et al. [11] performed a longitudinal study of manual and automatic facial expressions during semi-structured clinical interviews of 34 clinically depressed patients. They found that for both manual and automatic facial muscle activity analysis, participants with high symptom severity produced more expressions associated with contempt, smile less, and the smiles that were made were more likely to be related to contempt. Yang et al [29] analysed the vocal prosody of 57 participants of the same study. They found moderate predictability of the depression scores based on a combination of F_0 and switching pauses. Both studies used the Hamilton Rating Scale for Depression, which is a multiple choice questionnaire filled in by a clinician and used to provide an indication of depression, and as a guide to evaluate recovery. Scherer et al. [21] studied the correlation between automatic gaze, head pose, and smile detection and three mental health conditions (Depression, Post-Traumatic Stress Disorder and Anxiety). Splitting 111 participants into three groups based on their self-reported distress, they found significant differences for the automatically detected behavioural descriptors between the highest and lowest distressed groups.

Dimensional affect recognition aims to improve the understanding of human affect by modelling affect as a small number of continuously valued, continuous time signals. Compared to the more limited categorical emotion description (e.g. six basic emotions) and the computationally intractable appraisal theory, dimensional affect modelling has the benefit of being able to: a. encode small changes in affect over time, and b. distinguish between many more subtly different displays of affect, while remaining within the reach of current signal processing and machine learning capabilities.

The 2013 Audio-Visual Emotion Challenge and Workshop (AVEC 2013) will be the third competition event aimed at comparison of multimedia processing and machine learning methods for automatic audio, video and audiovisual emotion analysis, with all participants competing under strictly the same conditions. The goal of the AVEC Challenges series is to provide a common benchmark test set for individual multimedia processing and to bring together the audio and video emotion recognition communities, to compare the rel-

ative merits of the two approaches to emotion recognition under well-defined and strictly comparable conditions and establish to what extent fusion of the approaches is possible and beneficial. In addition, AVEC 2013 has the goal to accelerate the development of technologies that can aid the mental health profession in their aim to help people with mood disorders.

Following up from AVEC 2011 [25] and AVEC 2012 [24], which respectively used a categorical description of affect and automatic continuous affect recognition from audio and video on the SEMAINE database of natural dyadic interactions [18], we now aim to extend the analysis of affective behaviour to infer a more complex mental state, to wit, depression. Both the dimensional affect and the depression recognition problems are posed as a regression problem, and can thus be considered to be both challenging and rewarding. A major difference between this AVEC and the previous two is that the first two had as task making predictions on a very short temporal scale (either for every video frame or per spoken word), whereas AVEC 2013 extends this to include event recognition in the form of inferring a measure of depression for every recording.

Different from the continuous dimensional affect prediction, event-based recognition provides a single label over some pre-defined period of time rather than at every moment in time. In essence, continuous prediction is used for relatively fast-changing variables such as valence or arousal, while event-based recognition is more suitable for slowly varying variables such as mood or level of depression. One important aspect is that agreement must exist on what constitutes an event in terms of a logical unit in time. In this challenge, an event is defined as a participant performing a single experiment from beginning to end.

We are calling for teams to participate in emotion and depression recognition from video analysis, acoustic audio analysis, linguistic audio analysis, or any combination of these. As benchmarking database the DEPRESSION database of naturalistic video and audio of participants partaking in a human-computer interaction experiment will be used, which contains labels for the two target affect dimensions arousal and valence, and Beck Depression Index (BDI), a self-reported 21 multiple choice inventory [3].

Two Sub-Challenges are addressed in AVEC 2013:

- The *Affect Recognition Sub-Challenge (ASC)* involves fully continuous affect recognition of the dimensions valence and arousal (VA), where the level of affect has to be predicted for every moment of the recording.
- The *Depression Recognition Sub-Challenge (DSC)* requires participants to predict the level of self-reported depression as indicated by the BDI for every experiment session, that is, one continuous value per multimedia file.

For the ASC, two regression problems need to be solved for Challenge participation: prediction of the continuous dimensions AROUSAL and VALENCE. The ASC competition measure is the Pearson's correlation coefficient averaged over all sessions and both dimensions. For the DSC, a single regression problem needs to be solved. The DSC competition measure is root mean square error over all sessions.

Both Sub-Challenges allow contributors to find their own features to use with their regression algorithm. In addition,

standard feature sets are provided (for audio and video separately), which participants are free to use. The labels of the test partition remain unknown to the participants, and participants have to stick to the definition of training, development, and test partition. They may freely report on results obtained on the development partition, but are limited to five trials per Sub-Challenge in submitting their results on the test partition.

To be eligible to participate in the challenge, every entry has to be accompanied by a paper presenting the results and the methods that created them, which will undergo peer-review. Only contributions with an accepted paper will be eligible for Challenge participation. The organisers preserve the right to re-evaluate the findings, but will not participate in the Challenge themselves.

We next introduce the Challenge corpus (Sec. 2) and labels (Sec. 3), then audio and visual baseline features (Sec. 4), and baseline results (Sec. 5), before concluding in Sec. 6.

2. DEPRESSION DATABASE

The challenge uses a subset of the audio-visual depressive language corpus (AViD-Corpus), which includes 340 video clips of subjects performing a Human-Computer Interaction task while being recorded by a webcam and a microphone. There is only one person in every clip and the total number of subjects is 292, i.e. some subjects feature in more than one clip. The speakers were recorded between one and four times, with a period of two weeks between the measurements. Five subjects appear in four recordings, 93 in 3, 66 in 2, and 128 in only one sessions. The length of the clips is between 50 minutes and 20 minutes (mean = 25 minutes). The total duration of all clips is 240 hours. The mean age of subjects was 31.5 years, with a standard deviation of 12.3 years and a range of 18 to 63 years. The recordings took place in a number of quiet settings.

The behaviour within the clips consisted of different tasks which were Power Point guided: i.e., sustained vowel phonation, sustained loud vowel phonation, and sustained smiling vowel phonation; speaking out loud while solving a task; Counting from 1 to 10; reading out loud: excerpts of the novel “Homo” Faber by Max Frisch and the fable “Die Sonne und der Wind” (The North Wind and the Sun); singing: a German nursery rhyme “Guten Abend, gute Nacht” and “Aber bitte mit Sahne” by Udo Jürgens; telling a story from the subject’s own past: best present ever and sad event in the childhood; Telling an imagined story applying the Thematic Apperception Test (TAT), containing e.g. pictures of a man and a woman in bed, or a housewife and children who are trying to reach the cookies.

The audio was recorded using a headset connected to the built-in sound card of a laptop, at a sampling rate of 41 kHz, 16 bit. The original video was recorded using a variety of codecs and frame rates, and was resampled to a uniform 30 frames per second at 640 × 480 pixels, with 24 bits per pixels. The codec used was H.264, and the videos were embedded in an mp4 container.

For the organisation of the challenge, the recordings were split into three partitions: a training, development, and test set of 50 recordings each. The audio and audio-visual source files and the baseline features (see section 4) can be downloaded for all three partitions, but the labels are available only for the training and development partitions. All data

can be downloaded from a special user-level access controlled website (<http://avec2013-db.sspnet.eu>).

3. CHALLENGE LABELS

The affective dimensions used in the challenge were selected based on their relevance to the task of depression estimation. These are the dimensions AROUSAL and VALENCE, which form a well-established basis for emotion analysis in the psychological literature [10].

AROUSAL (Activity) is the individual’s global feeling of dynamism or lethargy. It subsumes mental activity, and physical preparedness to act as well as overt activity. VALENCE is an individual’s overall sense of “weal or woe”: Does it appear that, on balance, the person rated feels positive or negative about the things, people, or situations at the focus of his/her emotional state?

A team of 23 naive raters annotated all human-computer interactions. The raters annotated the two dimensions in continuous time and continuous value using a tool developed especially for this task. The annotations are often called *traces* after the early popular system that performed a similar function called FeelTrace [5]. An instantaneous annotation value is obtained using a two-axis joystick.

Every video was annotated by only a single rater for every dimension, due to time constraints. To reduce the annotators’ cognitive load (and hence improve annotation accuracy) each dimension (valence and arousal) was annotated separately. The annotation process resulted in a set of trace vectors $\{\mathbf{v}_i^a, \mathbf{v}_i^v\} \in \mathbb{R}$ for every rater i and dimension a (AROUSAL) and v (VALENCE). Every annotator was made to annotate a common reference video, which can be used to construct models that can compensate for inter-annotator variability in the remaining (singly-annotated) traces.

Sample values are obtained by polling the joystick in a tight loop. As such, inter-sample spacing is irregular (though minute, as implementation is in C++). These original traces are binned in temporal units of the same duration as a single video frame (i.e., 1/30 seconds). The raw joystick data for AROUSAL, and VALENCE lies in the range $[-1000, 1000]$ labels, which is normalised to the range $[-1, 1]$. The annotation tool used will be made available in the near-future.

The level of depression is labelled with a single value per recording using a standardised self-assessed subjective depression questionnaire, the Beck Depression Inventory-II (BDI-II, [3]). BDI-II contains 21 questions, where each is a forced-choice question scored on a discrete scale with values ranging from 0 to 3. Some items on the BDI-II have more than one statement marked with the same score. For instance, there are two responses under the Mood heading that score a 2: (2a) I am blue or sad all the time and I can’t snap out of it and (2b) I am so sad or unhappy that it is very painful. The final BDI-II scores range from 0 – 63: 0–13: indicates minimal depression, 14–19: indicates mild depression, 20–28: indicates moderate depression, 29–63: indicates severe depression.

The average BDI-level in the AVEC 2013 partitions was 15.1 for the training partition and 14.8 for the development partition (standard deviations = 12.3 and 11.8, respectively). For every recording in the training and development partitions a separate file with a single value is provided for the DSC, together with two files containing the affective dimension labels.

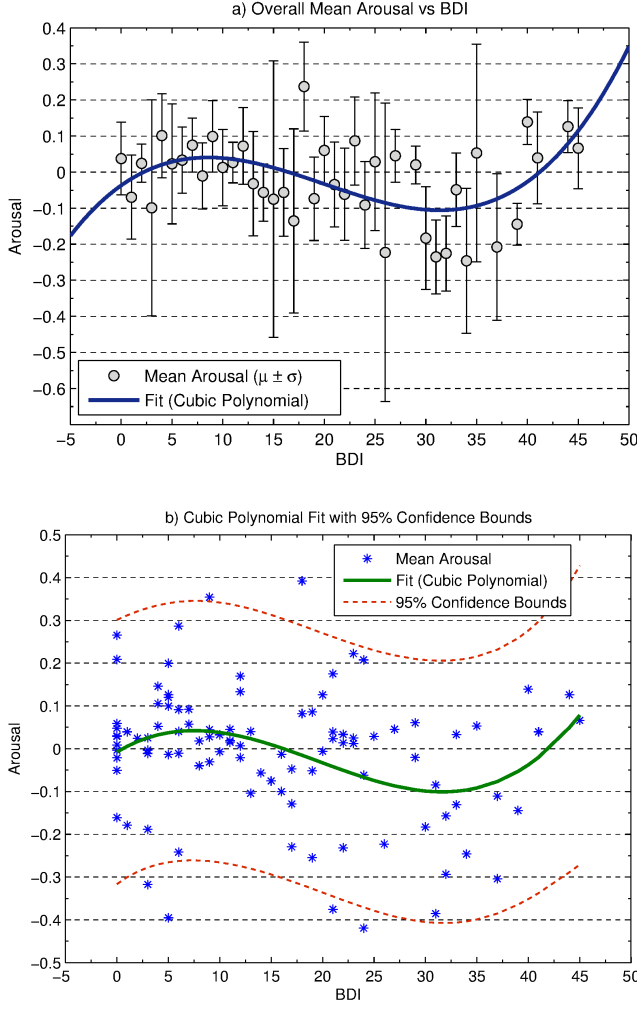


Figure 1: Cubic Polynomial correlation between Mean Arousal and BDI, demonstrated across a) Dataset mean values and b) Full raw dataset (with 95% confidence intervals).

To evaluate how strong the expected correlation between depression and the affective dimensions is, we calculated the average (mean) arousal and valence values over the duration of each recording in the training and development partitions. These values were then compared against their respective BDI ratings.

We observed a non-linear correlation between the depression and affect labels. Figures 1 and 2 show all data points for the training and development sets, where each data point is the mean arousal or valence over the whole video. The figures also show a 3 and a 5 degree polynomial for mean arousal respectively valence fit to BDI. Note what appears to be outliers in the mean valence for high BDI. While this seems to contradict the theory that depressed people have low valence, it is not uncommon for people with a high depression to display expressions of high valence.

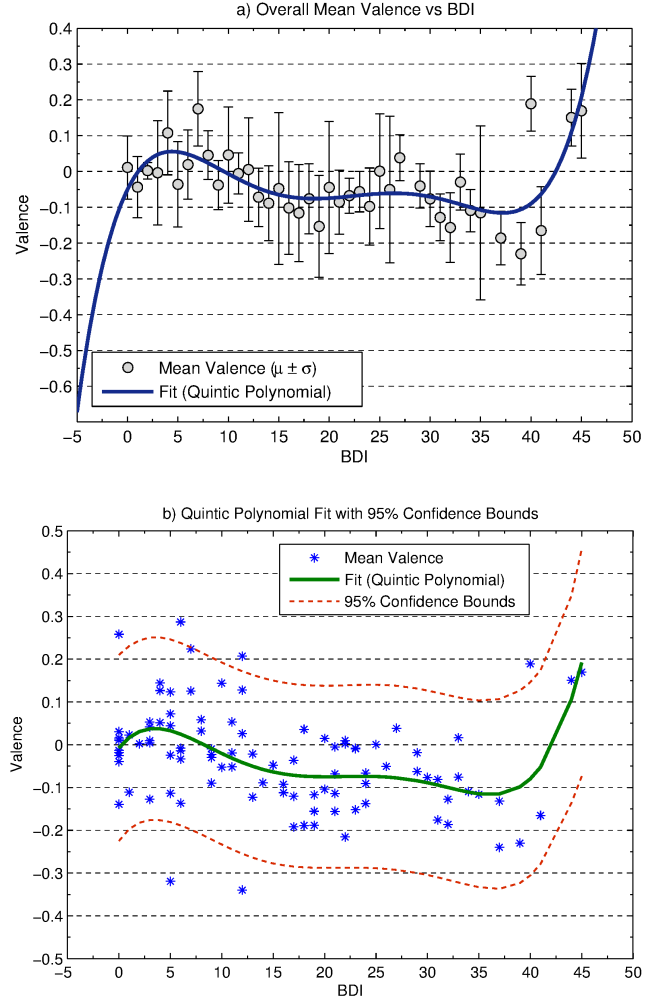


Figure 2: Quintic Polynomial correlation between Mean Valence and BDI, demonstrated across a) Dataset mean values and b) Full raw dataset (with 95% confidence intervals).

4. BASELINE FEATURES

In the following sections we describe how the publicly available baseline feature sets are computed for either the audio or the video data. Participants could use these feature sets exclusively or in addition to their own features.

4.1 Video Features

The majority of the features extracted from the video streams are dense local appearance descriptions. The descriptors that generate these features are most effective if they are applied to frontal faces of uniform size. Since the head pose and distance to the camera vary over time in the AVID-CORPUS recordings, we first detect the location of the face, and within that the locations of the eyes to help reduce the pose variance. The information describing the position and pose of the face and eyes are in themselves valuable for recognising the dimensional affect and are thus included with the set of video features together with the appearance descriptors. Fig. 3 gives an overview of the video feature extraction processing steps.

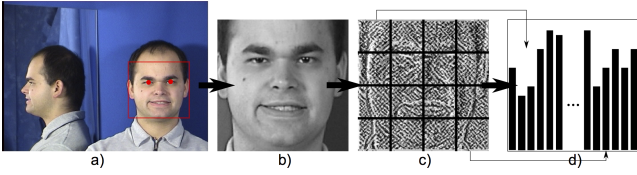


Figure 3: Video feature extraction overview: a) detection of face and eyes b) face normalised based on eye locations, c) extraction of LPQ features, d) divided in 4×4 blocks from which histograms of separate blocks are computed and concatenated into a single histogram.

To obtain the face position, we employ the open-source implementation of the Viola & Jones face detector that is included in OpenCV. This returns a four-valued descriptor of the face position and size. To wit, it provides the position of the top-left corner of the detected face area (f_x, f_y), followed by its width f_w and height f_h . The height and width output of this detector is rather unstable: Even in a video in which a face hardly moves the values for the height and width vary significantly (approximately 5% standard deviation). The face detector also doesn't provide any information about the head pose.

To refine the detected face region, and allow the appearance descriptor to correlate better with the shown expression rather than variations in head pose and face detector output, we proceed with detection of the locations of the eyes. This is again done with the OpenCV implementation of a Haar-cascade object detector, trained for either a left or a right eye. Let us define the detected left and right eye locations as p_l respectively p_r , and the line connecting p_l and p_r as l_e . The angle between l_e and the horizontal is then defined as α . The registered image is now obtained by rotating it so that $\alpha = 0$ degrees, then scaled to make the distance between the eye locations $\|p_l - p_r\| = 100$ pixels, and finally cropped to be 200 by 200 pixels, with p_r at position $\{p_r^x, p_r^y\} = \{80, 60\}$ to obtain the registered face image.

In AVEC 2011 and 2012, uniform Local Binary Patterns [19] were used as dense local appearance descriptors. For AVEC 2013, we chose instead to use Local Phase Quantisation (LPQ) as that was found to attain higher performance in facial expression recognition tasks [17]. The dynamic appearance descriptor LPQ-TOP was found to be even more accurate, but that descriptor depends on a near-perfect alignment of faces in subsequent frames, which is not possible in a near-real time automatic fashion on the AViD-Corpus dataset.

LPQs have been used extensively for face analysis in recent years, e.g., for face recognition [1], emotion detection [26], or detection of facial muscle actions (FACS Action Units) [17]. The LPQ descriptor extracts local phase information using the 2-D DFT or, more precisely, a short-term Fourier transform (STFT) computed over a rectangular M -by- M neighbourhood N_x at each pixel position \mathbf{x} of the image $f(\mathbf{x})$ defined by

$$F(\mathbf{u}, \mathbf{x}) = \sum_{\mathbf{y} \in N_x} f(\mathbf{x} - \mathbf{y}) e^{-j2\pi \mathbf{u}^T \mathbf{y}} = \mathbf{w}_{\mathbf{u}}^T \mathbf{f}_{\mathbf{x}} \quad (1)$$

where $\mathbf{w}_{\mathbf{u}}$ is the basis vector of the 2-D DFT at frequency \mathbf{u} , and $\mathbf{f}_{\mathbf{x}}$ is the vector containing all M^2 samples from N_x .

The local Fourier coefficients are computed at four frequency points: $u_1 = [a, 0]^T$, $u_2 = [0, a]^T$, $u_3 = [a, a]^T$, and $u_4 = [a, -a]^T$, where a is a sufficiently small scalar ($a = 1/M$ in our experiments). For each pixel position this results in a vector $F_x = [F(u_1, x), F(u_2, x), F(u_3, x), F(u_4, x)]$. The phase information in the Fourier coefficients is recorded by examining the signs of the real and imaginary parts of each component in F_x . This is done by using a simple scalar quantiser

$$q_j = \begin{cases} 1 & \text{if } g_j \geq 0 \text{ is true} \\ 0 & \text{otherwise} \end{cases} \quad (2)$$

where $g_j(x)$ is the j th component of the vector G_x and $G_x = [\text{Re}\{F_x\}, \text{Im}\{F_x\}]$. The resulting eight bit binary coefficients $q_j(x)$ are represented as integers using binary coding:

$$f_{\text{LPQ}}(x) = \sum_{j=1}^8 q_j 2^{j-1}. \quad (3)$$

As a result, a histogram of these values from all positions is composed to form a 256-dimensional feature vector. Histograms discard all information regarding the spatial arrangement of the patterns. In order to preserve some of this information, we divide the face region into 4×4 local regions, from which LPQ histograms are extracted and then concatenated into a single feature histogram (see Fig. 3). In the case of no face/eyes detected, the corresponding feature vector is set to all zeros so that it could be excluded in the training process or hold the last value when doing prediction.

4.2 Audio Features

In this Challenge, as was the case for AVEC 2012 and AVEC 2011, an extended set of features with respect to the INTERSPEECH 2009 Emotion Challenge (384 features) [22] and INTERSPEECH 2010 Paralinguistic Challenge (1582 features) [23] is given to the participants, again using the freely available open-source Emotion and Affect Recognition (openEAR) [7] toolkit's feature extraction backend openSMILE [8]. In contrast to AVEC 2011, the AVEC 2012 feature set was reduced by 100 features that were found to carry very little information, as they were zero or close to zero most of the time. In the AVEC 2013 feature set bugs in the extraction of jitter and shimmer were corrected, the spectral flatness was added to the set of spectral low-level descriptors (LLDs) and the MFCCs 11–16 were included in the set.

Thus, the AVEC 2013 audio baseline feature set consists of 2268 features, composed of 32 energy and spectral related low-level descriptors (LLD) \times 42 functionals, 6 voicing related LLD \times 32 functionals, 32 delta coefficients of the energy/spectral LLD \times 19 functionals, 6 delta coefficients of the voicing related LLD \times 19 functionals, and 10 voiced/unvoiced durational features. Details for the LLD and functionals are given in tables 1 and 2 respectively. The set of LLD covers a standard range of commonly used features in audio signal analysis and emotion recognition.

The audio features are computed on short episodes of audio data. As the data in the Challenge contains long continuous recordings, a segmentation of the data had to be performed. A set of baseline features is provided for three different versions of segmentation: First, a voice activity detector [6] was applied to obtain a segmentation based on speech activity. Pauses of more than 200 ms are used to

Table 1: 32 low-level descriptors.

Energy & spectral (32)
loudness (auditory model based), zero crossing rate, energy in bands from 250–650 Hz, 1 kHz–4 kHz, 25 %, 50 %, 75 %, and 90 % spectral roll-off points, spectral flux, entropy, variance, skewness, kurtosis, psychoacoustic sharpness, harmonicity, flatness, MFCC 1-16
Voicing related (6)
F_0 (sub-harmonic summation, followed by Viterbi smoothing), probability of voicing, jitter, shimmer (local), jitter (delta: “jitter of jitter”), logarithmic Harmonics-to-Noise Ratio (logHNR)

Table 2: Set of all 42 functionals. ¹Not applied to delta coefficient contours. ²For delta coefficients the mean of only positive values is applied, otherwise the arithmetic mean is applied. ³Not applied to voicing related LLD.

Statistical functionals (23)
(positive ²) arithmetic mean, root quadratic mean, standard deviation, flatness, skewness, kurtosis, quartiles, inter-quartile ranges, 1 %, 99 % percentile, percentile range 1 %–99 %, percentage of frames contour is above: minimum + 25%, 50%, and 90 % of the range, percentage of frames contour is rising, maximum, mean, minimum segment length ^{1,3} , standard deviation of segment length ^{1,3}
Regression functionals¹ (4)
linear regression slope, and corresponding approximation error (linear), quadratic regression coefficient a , and approximation error (linear)
Local minima/maxima related functionals¹ (9)
mean and standard deviation of rising and falling slopes (minimum to maximum), mean and standard deviation of inter maxima distances, amplitude mean of maxima, amplitude range of minima, amplitude range of maxima
Other^{1,3} (6)
LP gain, LPC 1–5

Table 3: Baseline results for affect recognition. Performance is measured in Pearson’s correlation coefficient averaged over all sequences.

Partition	Modality	Valence	Arousal	Average
Development	Audio	0.338	0.257	0.298
Development	Video	0.337	0.157	0.247
Test	Audio	0.089	0.090	0.089
Test	Video	0.076	0.134	0.105

split speech activity segments. Functionals are then computed over each detected segment of speech activity. These features can be used both for the emotion and depression tasks. The second segmentation method considers overlapping short fixed length segments (3 seconds) which are shifted forward at a rate of one second. These features are intended for the emotion task. The third method also uses overlapping fixed length segments shifted forward at a rate of one second, however, the windows are 20 seconds long to capture slow changing, long range characteristics. These features are expected to perform best in the depression task.

5. CHALLENGE BASELINES

For transparency and reproducibility, we use standard algorithms. We conducted two separate baselines: one using video features only, and the other using only audio features.

For the video-based baseline, a combination of geometric features (head location, head motion, head pose) and appearance features (LPQ) is employed. The geometric features are the head location (with respect to the first frame), head motion (with respect the previous frame), head pose (roll), and head pose changes (with respect to the previous frame). The geometric features are computed using the face and eye detection results included in the baseline features. As appearance features we used the LPQ features included in the baseline features.

To deal with variability in appearance and registration errors, at each selected training example we used a window of 100 frames, taking the mean of the extracted features over these frames as the feature vector. To deal with the large number of frames in the training set, we down sampled the number of training features by a factor 4. The baseline method uses correlation based feature selection (CFS) [12] and ϵ -Support Vector Machine Regressors (ϵ -SVR) with an intersection kernel. This means only two variables need to be optimised, that is, the slack variable C and ϵ itself. For results on the development set, optimisation is done in a 5-fold cross-validation loop on the training data. The results of the affect recognition are shown in Table 3.

As the DSC has a single label per video, we created a single feature vector per video by taking the median value of all video features. Again, CFS was used to reduce the number of features, in conjunction with SVRs with intersection kernels. Results for video are shown in Table 4. Please note that participants in the challenge will be ranked based on their RMSE results. We compared the baseline predictions to a naive or chance level error. This was obtained by calculating the average BDI value over either the training set (for prediction on the development partition) or over both the training and development sets (for prediction on the test partition), and using that as the predicted depression level

Table 4: Baseline results for depression recognition. Performance is measured in mean absolute error (MAE) and root mean square error (RMSE) over all sequences.

Partition	Modality	MAE	RMSE
Development	Audio	8.66	10.75
Development	Video	8.74	10.72
Test	Audio	10.35	14.12
Test	Video	10.88	13.61

for all sessions we test on. For depression prediction on the development partition, chance levels would be an error of 11.90 RMSE and 10.28 MAE.

For the audio-based baseline, features extracted from short 3-second segments performed best for valence, while for arousal features from automatically detected voice activity [6] segments worked better. For audio, SVRs with a linear kernel were used. For depression recognition best results were obtained by computing the audio features over longer 20 second non-overlapping segments, which were subsequently averaged over the entire recording. Results are shown in Tables 3 and 4.

6. CONCLUSION

We introduced AVEC 2013 – the first combined open Audio/Visual Emotion and Depression recognition Challenge. It addresses in two sub-challenges the detection of the affective dimensions valence and arousal in continuous time and value, and the estimation of a self-reported level of depression. This manuscript describes AVEC 2013’s challenge conditions, data, baseline features and results. By intention, we opted to use open-source software and the highest possible transparency and realism for the baselines by re-training from feature space optimisation and optimising on test data. This should improve the reproducibility of the baseline results.

Acknowledgments

The authors would like to thank Prof Jeff Cohn for his feedback on our draft manuscript, and Sander Koelstra for his invaluable technical support in setting up the web-accessible challenge database. The work of Michel Valstar is partly funded by the NIHR-HTC ‘MindTech’ and Horizon Digital Economy Research, RCUK grant EP/G065802/1. The work of Sebastian Schnieder is partly funded by the German Research Foundation (KR3698/4-1). The challenge in general has been generously supported by the HUMAINE organisation and the EU network of excellence on Social Signal Processing SSPNet (European Community’s 7th Framework Programme [FP7/20072013] under grant agreement no. 231287). The authors further acknowledge funding from the European Commission(grant no. 289021, ASC-Inclusion).

7. REFERENCES

- [1] T. Ahonen, A. Hadid, and M. Pietikäinen. Face description with local binary patterns: Application to face recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 28(12):2037–2041, 2006.
- [2] M. R. Bagby, A. G. Ryder, D. R. Schuller, and M. B. Marshall. The hamilton depression rating scale: Has the gold standard become a lead weight? *American Journal of Psychiatry*, 161:2163–2177, 2004.
- [3] A. Beck, R. Steer, R. Ball, and W. Ranieri. Comparison of beck depression inventories -ia and -ii in psychiatric outpatients. *Journal of Personality Assessment*, 67(3):588–97, December 1996.
- [4] J. F. Cohn, S. Kreuz, I. Matthews, Y. Yang, M. H. Nguyen, M. Tejera Padilla, and et al. Detecting depression from facial actions and vocal prosody. In *Proc. Affective Computing and Intelligent Interaction*, pages 1–7, 2009.
- [5] R. Cowie, E. Douglas-Cowie, S. Savvidou, E. McMahon, M. Sawey, and M. Schröder. Feeltrace: An instrument for recording perceived emotion in real time. In *Proc. ISCA Workshop on Speech and Emotion*, pages 19–24, Belfast, UK, 2000.
- [6] F. Eyben, F. Wenginger, S. Squartini, and B. Schuller. Real-life voice activity detection with lstm recurrent neural networks and an application to hollywood movies. In *Proc. of ICASSP, Vancouver, Canada*. IEEE, 2013. to appear.
- [7] F. Eyben, M. Wöllmer, and B. Schuller. openEAR - Introducing the Munich Open-Source Emotion and Affect Recognition Toolkit. In *Proc. ACII*, pages 576–581, Amsterdam, The Netherlands, 2009.
- [8] F. Eyben, M. Wöllmer, and B. Schuller. openSMILE – The Munich Versatile and Fast Open-Source Audio Feature Extractor. In *Proc. ACM Multimedia (MM)*, pages 1459–1462, Florence, Italy, 2010.
- [9] M. First, R. Spitzer, M. Gibbon, and J. Williams. *Structured Clinical Interview for DSM-IV Axis I Disorders SCID-I: Clinician Version, Administration Booklet*. SCID-I: Clinician Version. American Psychiatric Press, 1997.
- [10] J. Fontaine, S. K.R., E. Roesch, and P. Ellsworth. The world of emotions is not two-dimensional. *Psychological science*, 18(2):1050 – 1057, 2007.
- [11] J. M. Girard, J. F. Cohn, M. H. Mahoor, S. Mavadati, and D. Rosenwald. Social risk and depression: Evidence from manual and automatic facial expression analysis. In *IEEE International Conference on Automatic Face and Gesture Recognition*, 2013.
- [12] M. Hall. *Correlation-based Feature Selection for Machine Learning*. PhD thesis, The University of Waikato, 1999.
- [13] M. Hamilton. Development of a rating scale for primary depressive illness. *British Journal of Social and Clinical Psychology*, 8:278–296, 1967.
- [14] J. Hamm, Kohler, C. G., Gur, R. C., and R. Verma. Automated facial action coding system for dynamic analysis of facial expressions in neuropsychiatric disorders. *Journal of Neuroscience Methods*, 200(2):237–256, 2011.
- [15] Health & Consumer Protection Directorate General. Improving the mental health of the population: Towards a strategy on mental health for the european union. Technical report, European Union, 2005.
- [16] Health & Consumer Protection Directorate General. Mental health in the eu. Technical report, European Union, 2008.

- [17] B. Jiang, M. Valstar, and M. Pantic. Action unit detection using sparse appearance descriptors in space-time video volumes. In *Proc. IEEE Int. Conf. on Automatic Face and Gesture Recognition*, pages 314–321, Santa Barbara, USA, 2011.
- [18] G. McKeown, M. Valstar, R. Cowie, M. Pantic, and M. Schroder. The semaine database: Annotated multimodal records of emotionally colored conversations between a person and a limited agent. *IEEE Trans. Affective Computing*, 3:5–17, 2012.
- [19] T. Ojala, M. Pietikainen, and T. Maenpaa. Multiresolution gray-scale and rotation invariant texture classification with local binary patterns. *IEEE Trans. Pattern Analysis and Machine Intelligence*, 24(7):971–987, 2002.
- [20] R. Picard. *Affective Computing*. MIT Press, 1997.
- [21] S. Scherer, G. Stratou, J. Gratch, J. Boberg, M. Mahmoud, A. S. Rizzo, and L.-P. Morency. Automatic behavior descriptors for psychological disorder analysis. In *Proc. IEEE Int. Conf. on Automatic Face and Gesture Recognition*, 2013.
- [22] B. Schuller, S. Steidl, and A. Batliner. The INTERSPEECH 2009 Emotion Challenge. In *Proc. INTERSPEECH 2009*, pages 312–315, Brighton, UK, 2009.
- [23] B. Schuller, S. Steidl, A. Batliner, F. Burkhardt, L. Devillers, C. Müller, and S. Narayanan. The INTERSPEECH 2010 Paralinguistic Challenge. In *Proc. INTERSPEECH 2010*, pages 2794–2797, Makuhari, Japan, 2010.
- [24] B. Schuller, M. Valstar, F. Eyben, R. Cowie, and M. Pantic. AVEC 2012 - the continuous audio/visual emotion challenge. In *Proceedings ACM Int'l Conf. Multimodal Interaction*, pages 449–456, October 2012.
- [25] B. Schuller, M. Valstar, F. Eyben, G. McKeown, R. Cowie, and M. Pantic. AVEC 2011 - The First International Audio/Visual Emotion Challenge. In *Proceedings Int'l Conference on Affective Computing and Intelligent Interaction 2011, ACII 2011*, volume II, pages 415–424, Memphis, TN, October 2011. Springer.
- [26] C. Shan, S. Gong, and P. W. Mcowan. Facial expression recognition based on local binary patterns: A comprehensive study. *Image and Vision Computing*, 27(6):803–816, 2009.
- [27] A. Vinciarelli, M. Pantic, D. Heylen, C. Pelachaud, I. Poggi, F. D'ericco, and M. Schroeder. Bridging the gap between social animal and unsocial machine: A survey of social signal processing. *IEEE Trans. Affective Computing*, 3:69–87, April 2012. Issue 1.
- [28] P. Wang, F. Barrett, E. Martin, M. Milonova, R. E. Gur, R. C. Gur, C. Kohler, and et al. Automated video-based facial expression analysis of neuropsychiatric disorders. *Journal of Neuroscience Methods*, 168(1):224 – 238, 2008.
- [29] Y. Yang, C. Fairbairn, and J. Cohn. Detecting depression severity from intra- and interpersonal vocal prosody. *IEEE Transactions on Affective Computing*, 4, 2013.
- [30] M. Zimmerman, I. Chelminski, and M. Posternak. A review of studies of the hamilton depression rating scale in healthy controls: Implications for the definition of remission in treatment studies of depression. *Journal of Nervous & Mental Disease*, 192(9):595–601, 2004.