

AUTOMATIC DETECTION OF SENTIMENTALITY FROM FACIAL EXPRESSIONS

Mina Bishay, Jay Turcot, Graham Page and Mohammad Mavadati

Smart Eye AB

ABSTRACT

Emotion recognition has received considerable attention from the Computer Vision community in the last 20 years. However, most of the research focused on analyzing the six basic emotions (e.g. joy, anger, surprise), with a limited work directed to other affective states. In this paper, we tackle sentimentality (strong feeling of heartwarming or nostalgia), a new emotional state that has few works in the literature, and no guideline defining its facial markers. To this end, we first collect a dataset of 4.9K videos of participants watching some sentimental and non-sentimental ads, and then we label the moments evoking sentimentality in the ads. Second, we use the ad-level labels and the facial Action Units (AUs) activation across different frames for defining some weak frame-level sentimentality labels. Third, we train a Multilayer Perceptron (MLP) using the AUs activation for sentimentality detection. Finally, we define two new ad-level metrics for evaluating our model performance. Quantitative and qualitative results show promising results for sentimentality detection. To the best of our knowledge this is the first work to address the problem of sentimentality detection.

Index Terms— Sentimentality, Facial expressions, AU detection, Ad-level KPIs.

1. INTRODUCTION

Understanding facial expressions are quite important to analyze humans’ emotions and non-verbal communications. Automatic Facial Expression Analysis (AFEA) has been an active research area in Computer Vision, as it has gained popularity in several applications like ad testing [1, 2, 3], driver state monitoring [4, 5], and health care [6, 7, 8]). In ad testing, analyzing customers’ facial responses gives traders insights about customers engagement, liking, and purchase intent [1]. However, AFEA is quite limited to the detection of AUs and basic emotions, as collecting and labelling real-world data for other emotional states are quite challenging.

Evoking emotions like sentimentality (an emotion with heartwarming or nostalgic feelings) in commercial ads is an emerging trend in advertising [9]. Subsequently, there has been a growing interest in studying sentimentality [10, 11]. In [10], McDuff highlighted the prominent AUs in sentimental responses, while in [11] McDuff classified the ad

media content into 5 classes (informed, inspired, sentimental, amused, persuaded) based on the participants facial responses and some media features. The facial markers of sentimentality has not been defined in the literature, in contrast to the basic emotions that were interpreted from AUs through the Emotional Facial Action Coding System (EMFACS) [12]. Therefore, it was difficult to directly label and detect sentimentality.

In this paper, we present a novel methodology for detecting sentimentality. Specifically, we first collect real-world facial responses for participants watching a group of sentimental and non-sentimental ads. Second, we label the moments evoking sentimentality in the ads. Third, using the ad-level labels and an AU detector (predicting 20 different AUs), we filter and categorize the participants’ frames in the training set into positive and negative examples. Specifically, frames with active AUs shown during sentimental moments are considered positive sentimentality examples, while other frames negative examples. Finally, we use the frame-level labels and AU predictions for training a MLP, that extracts high-level features on the top of the AU predictions (low-level features) for sentimentality detection. Fig. 1 shows an overview of our architecture. We believe that the proposed methodology can be replicated to other untackled emotions without the need for the exhaustive frame-level labelling, nor predefined facial markers for the target emotion.

For evaluating our architecture, we define two new ad-level Key Performance Indicators (KPIs), that are based on the ad-level aggregated sentimentality. The first KPI measures how separable are the sentimental and non-sentimental ads in terms of the aggregated sentimentality, while the second measures if the aggregated sentimentality is firing high at the right sentimental moments. Our architecture shows promising qualitative and quantitative results for sentimentality detection. We base our emotion/sentimentality analysis on the recommendations given by Barrett *et al.* in [13] for studying facial movements in real life, sampling across different cultures, and using multiple facial stimuli.

To the best of our knowledge this is the first work to directly address the problem of sentimentality detection. It is worth noting that in this paper we detect “sentimentality”, a kind of emotion with heartwarming or nostalgic feelings, from an image capturing participant’s facial expression – this is completely different from “sentiment analysis” that has a

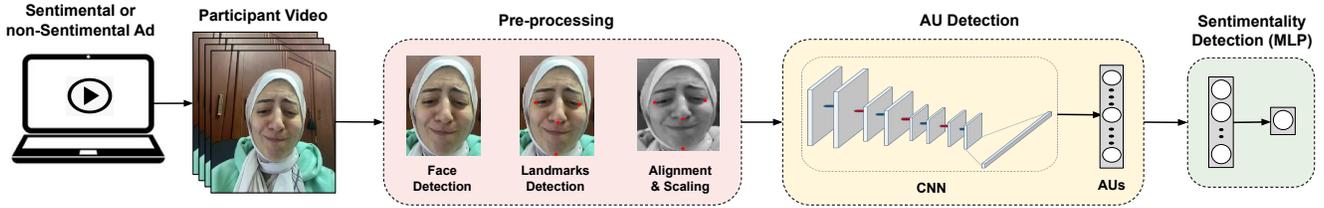


Fig. 1. The proposed architecture for sentimentality detection.

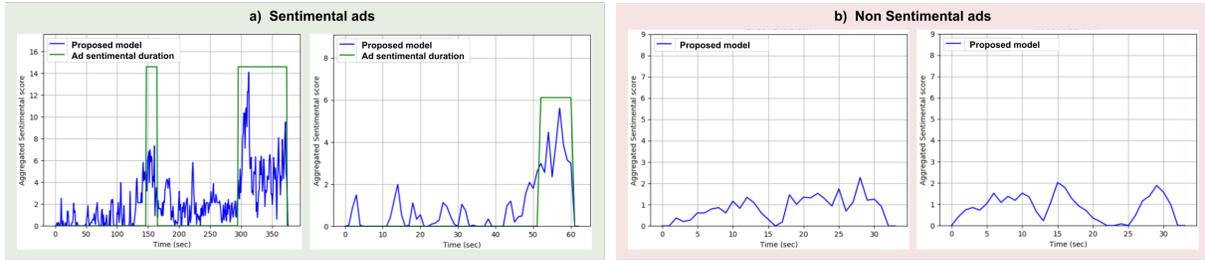


Fig. 2. The aggregated sentimentality across different sentimental and non-sentimental ads.

quite good amount of work in the literature [14, 15], and basically detects if a piece of writing (e.g. reviews, survey responses) was positive, negative, or neutral.

The rest of the paper is organized as follows: In Section 2 and Section 3, we present the sentimentality dataset and the ad-level KPIs used in our analysis, respectively. In Section 4, we introduce the proposed methodology for sentimentality detection. Finally, we draw our conclusions in Section 5.

2. SENTIMENTALITY DATASET

Affectiva/SmartEye in collaboration with global market agencies have collected and analyzed thousands of commercial ads across different markets. For each ad, participants were hired to watch the ad, and then fill a survey about how they feel about the ad. A consent was given by the participants to get video recorded while they were watching the ad. The participants’ facial responses in the videos were detected and analyzed to get insights about their level of engagement, liking, and purchase intent [1].

For our analysis, experienced ad testers have selected 33 ads (18 sentimental and 15 non-sentimental), and 4.9K participants’ videos to form a dataset for sentimentality. Sentimental ads are evoking sentimentality at some moments, while non-sentimental ads are typically informative, funny, or musical ads. The selected sentimental and non-sentimental ads span different markets (USA, UK, East and South Asia and Latin America), and subsequently the participants’ videos used in our analysis have diverse demographics (gender, age band and ethnicity). Note that different participants were recruited for watching the different ads, and those participants were not informed about the emotion the ad is trying to evoke.

The start and the end of the sentimental moments in the sentimental ads were labelled by 3 labellers, who were labelling by just watching the ads.

We use 3 sentimental ads for training our model, and 15 sentimental and 15 non-sentimental ads for testing. More ads/samples are needed in the testing as the proposed KPIs are calculated on the top of the aggregated sentimentality across different ads. For the participants’ videos, we have around 2.1K videos for the 18 sentimental ads (250 videos for the training ads and 1.85K videos for the testing ads), and 2.8K videos for the 15 non-sentimental ads.

3. AD-LEVEL KPIs

Our dataset does not have frame-level ground truth labels for sentimentality, so it is challenging to use the typical frame-level KPIs in our analysis. As the dataset has mainly ad-level labels defining the sentimental and non-sentimental ads, and the sentimental moments in the sentimental ads, we define new KPIs based on the available ad-level labels. Specifically, we aggregate (i.e. average) the participants’ predicted sentimentality across each sentimental and non-sentimental ad, in order to get a single sentimentality curve for each ad (Fig. 2 shows the aggregated sentimentality across four ads). Then, we calculate two KPIs that compare the ad-level predicted sentimentality to the ad-level labels.

The first KPI, named *ROC-Ad*, uses the area under the ROC curve (ROC-AUC) for measuring how separable are the sentimental and non-sentimental ads. To do so, we first calculate the maximum sentimentality score across the aggregate curve of each ad – this leads to 15 scores across sentimental ads (considered the positive predictions), and 15 scores

across non-sentimental ones (the negative predictions). Then, the ROC-AUC is calculated between the positive and negative predictions.

The second KPI, named *ROC-Sent*, measures if the model is firing high at the right sentimental moments. Specifically, ROC-Sent uses ROC-AUC for measuring how separable are the sentimental and non-sentimental moments in the sentimental ads. Similar to the ROC-Ad, we calculate the maximum sentimentality score across the 15 sentimental moments (positive predictions), and the 15 non-sentimental moments (negative predictions). Then, the ROC-AUC is calculated between the positive and negative predictions.

4. METHOD

In this section we present the proposed methodology for detecting sentimentality. Our model takes as input a frame depicting a face and gives as output a binary label indicating if the face is showing markers of sentimentality or not. The analysis is performed in 3 stages; preprocessing, AU detection, and sentimentality detection. Fig. 1 shows an overview of the whole architecture.

After collecting a dataset of diverse participants watching some sentimental and non-sentimental ads, and labeling the moments evoking sentimentality. We first detect and align the participant face. Then, we build an AU detector for analyzing the participant facial expression. The AU predictions are used for defining some weak frame-level sentimentality labels, as well as extracting low-level features from the face image. Finally, we train a MLP using the AU predictions for extracting high-level features for sentimentality detection.

4.1. Preprocessing

In order to process each video in our dataset, we first extract the region of interest (i.e. the participant face) at each frame, by using a face detector trained in the wild. To ensure we are including participants who are not distracted from the ad or away from the screen, we only include participants with face coverage (i.e. the percentage of the video frames with a face got detected) $\geq 90\%$ for the next steps. Second, we extract 4 facial landmarks (outer eye corners, nose tip and chin) from the face region. These landmarks are used for aligning the face horizontally to have a zero roll angle. Finally, the aligned faces are scaled to a fixed resolution, and passed as an input to the AU detection architecture.

4.2. AU Detection

In this section we describe the dataset, CNN architecture, and experimental settings used for building and training our AU detector. Most of the settings in our AU detection architecture are chosen based on the recommendations given in [17, 18].

AU Dataset. In the literature, several datasets (e.g. DISFA [19], UNBC [20]) have been used for training different architectures. However, many of these datasets have relatively limited number of participants, recording conditions, and/or diversity in demographics. In our analysis, we use a large-scale dataset consisting of $\sim 55K$ videos, that were captured in the wild. The participants in our dataset have diverse age, gender and ethnicity. For the experiments, we divide the dataset into 40.9K videos for training, 5.9K for validation, and 8.2K for testing.

Our large-scale dataset was collected using the web-based approach described in [21, 22]. The videos were collected worldwide (from 90+ countries) for participants watching commercial ads. The videos were manually annotated for the presence of 20 AUs using trained FACS coders. A part of this dataset was made available to the research community through AM-FED [21] and AM-FED+ [22]. Note that the videos in the AU dataset are different from the ones in the sentimentality dataset.

CNN architecture. We treat the AU detection problem as a multi-label classification problem where a single CNN is jointly trained for detecting 20 AUs simultaneously (AUs are given in Table 1). The CNN consists of 5 convolutional and 1 fully-connected layers. A max-pooling layer is used after each convolutional layer. The fully-connected layer has 20 sigmoid units representing the predictions of the 20 AUs.

Experimental settings. As we are using a naturalistic dataset, most of the AUs in our dataset are severely imbalanced, having a high ratio of negative to positive examples. In order to avoid biasing the classifier to the most frequent class, we use oversampling to balance the data. The training batches are augmented by flipping, shifting, rotation, etc. The Binary Cross Entropy function is used for calculating the loss.

We compare the developed AU detector to a widely-used facial analysis toolkit (AFFDEX-SDK [16]) on our large testing set. ROC-AUC is used for evaluating the AU detection performance. Table 1 shows the performance across the different AUs for both models. Our AU detector achieves better performance than AFFDEX-SDK for most of the AUs (by $\sim 4\%$ on average).

4.3. Sentimentality detection

We train a MLP on the top of the AU predictions for detecting sentimentality. For training the MLP, we need to have a set of positive and negative examples for sentimentality. There is no guideline in the literature that has defined the markers of sentimentality – this is in contrast to basic emotions that are described by AUs through EMFACS [12]. Subsequently, it is challenging to directly label a face dataset for sentimentality. In this paper we use the ad-level labels defining the sentimental moments for extracting positive and negative examples for sentimentality. Specifically, we use the expressive faces shown during sentimental moments as positive exam-

Table 1. Comparing the ROC-AUC of the developed AU detector to AFFDEX-SDK [16].

Facial Expressions	AU1	AU2	AU4	AU5	AU6	AU7	AU9	AU10	AU14	AU15	AU17	AU18	AU20	AU24	AU25	AU26	AU28	Eye closure	Smile	Smirk	Avg
AFFDEX [16]	0.76	0.79	0.86	0.87	0.92	0.75	0.91	0.86	0.86	0.78	0.79	0.91	0.86	0.76	0.86	0.63	0.91	0.92	0.94	0.82	0.84
Ours	0.79	0.84	0.92	0.85	0.94	0.83	0.89	0.93	0.80	0.88	0.88	0.92	0.93	0.87	0.89	0.71	0.96	0.91	0.97	0.84	0.88

Table 2. Sentimentality evaluation results across the different AUs and the proposed model.

KPIs	Chance level	AU1	AU2	AU4	AU5	AU6	AU7	AU9	AU10	AU14	AU15	AU17	AU18	AU20	AU24	AU25	AU26	AU28	Eye Closure	Smile	Smirk	Proposed
ROC-Ad	0.50	0.61	0.73	0.55	0.74	0.84	0.54	0.52	0.61	0.70	0.69	0.72	0.73	0.67	0.73	0.60	0.75	0.74	0.43	0.66	0.72	0.79
ROC-Sent	0.50	0.44	0.58	0.41	0.36	0.36	0.46	0.32	0.52	0.48	0.47	0.36	0.40	0.51	0.40	0.35	0.54	0.48	0.27	0.60	0.50	0.61
Avg	0.50	0.52	0.65	0.48	0.55	0.60	0.50	0.42	0.56	0.59	0.58	0.54	0.56	0.59	0.56	0.47	0.64	0.61	0.35	0.63	0.61	0.70

**Fig. 3.** The detected positive moments of sentimentality using the proposed model.

ples, and all the expressions shown during non-sentimental moments as negative examples.

In order to make the positive examples cleaner, we first discard all the frames with no active AUs. Second, we evaluate using the ad-level KPIs how indicative are the 20 detected AUs for sentimentality (i.e. test if the activation of a single AU can be used as a marker for sentimentality). Table 2 shows the results across the different AUs. Results show that most of the AUs are achieving relatively low performance on one or two of the KPIs. Subsequently, sentimentality is potentially expressed by more complex combination of AUs. Based on that, we discard all the positive frames with only one active AU. Eventually, positive examples include frames with ≥ 2 active AUs during sentimental moments, while the negative examples include any other frame (with or without active AUs) in the non-sentimental moments.

The MLP consists of 2 Fully-Connected (FC) layers, the first FC has 8 neurons for extracting high-level features on the top of the AU predictions, while the second has 1 neuron for detecting sentimentality. We train the MLP using the positive and negative examples extracted from the participants watching the 3 sentimental ads in the training set. We train the MLP for 100 epochs using the Adam optimizer. For testing, we use 15 sentimental and 15 non-sentimental ads. The detected sentimentality for different participants is aggregated for each ad in the testing set, and then the ad-level KPIs are calculated.

Table 2 shows the ROC-Ad and ROC-Sent achieved by the proposed model. On average our model has better performance than the chance level and the 20 AUs. The MLP combines different AUs to get better representation for sentimentality. Fig. 2 shows the aggregated sentimentality across 2 sentimental and 2 non-sentimental ads from the testing set, as well as the sentimental moments in the ads. The aggregate

curves show that our model has relatively higher activation at the sentimental ads than the non-sentimental ones. In addition, the model is firing more accurately at the sentimental moments. Fig. 3 shows some of the faces with positive sentimentality detection (faces belong to Smart Eye employees who have been recorded while watching a sentimental ad). Reviewing the detected faces shows that sentimentality has different combinations of AUs, and is not expressed in the same way across different people.

We believe that the proposed methodology can be replicated for other emotional states that has no predefined facial markers. In addition, collecting and labelling some stimuli (e.g. ads) evoking the target emotional state can replace the exhaustive frame-level labelling.

5. CONCLUSION

In this work we present a novel methodology for detecting sentimentality. Our architecture consists of 3 steps; a) detecting and aligning the participants' faces, b) detecting 20 different AUs for each frame (low-level features), and c) training a MLP for detecting sentimentality by extracting high-level features on the top of the AU predictions. For training our model, we first collect a dataset of participants watching some sentimental and non-sentimental ads, and then we label the moments evoking sentimentality in the ads. The ad-level labels along with the frame-level AU activations are used for defining a group of positive and negative examples for training the MLP. We define two new ad-level KPIs for evaluating our model performance, by measuring how separable are the sentimental and non-sentimental ads, and the sentimental and non-sentimental moments. Qualitative and quantitative results show promising results for sentimentality detection.

6. REFERENCES

- [1] Daniel McDuff, Rana El Kaliouby, et al., “Predicting ad liking and purchase intent: Large-scale analysis of facial responses to ads,” *IEEE Transactions on Affective Computing*, vol. 6, no. 3, pp. 223–235, 2014.
- [2] Daniel McDuff and Rana El Kaliouby, “Applications of automated facial coding in media measurement,” *IEEE transactions on affective computing*, vol. 8, no. 2, pp. 148–160, 2016.
- [3] Natalia Efremova, Navid Hajimirza, et al., “Understanding consumer attention on mobile devices,” in *FG 2020*. IEEE, 2020, pp. 919–919.
- [4] Alexandru Mălăescu, Liviu Cristian Duțu, et al., “Improving in-car emotion classification by nir database augmentation,” in *2019 14th IEEE International Conference on Automatic Face & Gesture Recognition (FG 2019)*. IEEE, 2019, pp. 1–5.
- [5] Torsten Wilhelm, “Towards facial expression analysis in a driver assistance system,” in *2019 14th IEEE International Conference on Automatic Face & Gesture Recognition (FG 2019)*. IEEE, 2019, pp. 1–4.
- [6] Shashank Jaiswal, Michel F Valstar, et al., “Automatic detection of adhd and asd from expressive behaviour in rgb data,” in *FG 2017*. IEEE, 2017, pp. 762–769.
- [7] Mina Bishay, Petar Palasek, Stefan Priebe, and Ioannis Patras, “Schinet: Automatic estimation of symptoms of schizophrenia from facial behaviour analysis,” *IEEE Transactions on Affective Computing*, 2019.
- [8] Mina Bishay, Stefan Priebe, and Ioannis Patras, “Can automatic facial expression analysis be used for treatment outcome estimation in schizophrenia?,” in *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2019, pp. 1632–1636.
- [9] Daniel McDuff, May Amr, and Rana El Kaliouby, “The fine line of sentimental advertising,” <https://blog.affectiva.com/the-fine-line-of-sentimental-advertising>, 2015.
- [10] Daniel McDuff, “Discovering facial expressions for states of amused, persuaded, informed, sentimental and inspired,” in *Proceedings of the International Conference on Multimodal Interaction*, 2016, pp. 71–75.
- [11] Daniel McDuff and Mohammad Soleymani, “Large-scale affective content analysis: Combining media content features and facial reactions,” in *2017 12th IEEE International Conference on Automatic Face & Gesture Recognition (FG 2017)*. IEEE, 2017, pp. 339–345.
- [12] Wallace V Friesen, Paul Ekman, et al., “Emfacs-7: Emotional facial action coding system,” *Unpublished manuscript, University of California at San Francisco*, vol. 2, no. 36, pp. 1, 1983.
- [13] Lisa Feldman Barrett, Ralph Adolphs, et al., “Emotional expressions reconsidered: Challenges to inferring emotion from human facial movements,” *Psychological science in the public interest*, vol. 20, no. 1, pp. 1–68, 2019.
- [14] Walaa Medhat, Ahmed Hassan, and Hoda Korashy, “Sentiment analysis algorithms and applications: A survey,” *Ain Shams engineering journal*, vol. 5, no. 4, pp. 1093–1113, 2014.
- [15] Lei Zhang, Shuai Wang, and Bing Liu, “Deep learning for sentiment analysis: A survey,” *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, vol. 8, no. 4, pp. e1253, 2018.
- [16] Daniel McDuff, Abdelrahman Mahmoud, et al., “Affdex sdk: a cross-platform real-time multi-face expression recognition toolkit,” in *Proceedings of the 2016 CHI conference extended abstracts on human factors in computing systems*, 2016, pp. 3723–3726.
- [17] Mina Bishay, Ahmed Ghoneim, et al., “Choose settings carefully: Comparing action unit detection at different settings using a large-scale dataset,” in *2021 IEEE International Conference on Image Processing (ICIP)*. IEEE, 2021, pp. 2883–2887.
- [18] Mina Bishay, Ahmed Ghoneim, et al., “Which cnns and training settings to choose for action unit detection? a study based on a large-scale dataset,” in *2021 16th IEEE International Conference on Automatic Face and Gesture Recognition (FG 2021)*. IEEE, 2021, pp. 1–5.
- [19] S Mohammad Mavadati, Mohammad H Mahoor, et al., “Disfa: A spontaneous facial action intensity database,” *IEEE Transactions on Affective Computing*, vol. 4, no. 2, pp. 151–160, 2013.
- [20] Patrick Lucey, Jeffrey F Cohn, et al., “Painful data: The unbc-mcmaster shoulder pain expression archive database,” in *FG 2011*. IEEE, 2011, pp. 57–64.
- [21] Daniel McDuff, Rana Kaliouby, et al., “Affectiva-mit facial expression dataset (am-fed): Naturalistic and spontaneous facial expressions collected,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2013, pp. 881–888.
- [22] Daniel McDuff, May Amr, and Rana El Kaliouby, “Am-fed+: An extended dataset of naturalistic facial expressions collected in everyday settings,” *IEEE Transactions on Affective Computing*, vol. 10, no. 1, pp. 7–17, 2018.