

Confusion Mental State Inference through Intel RealSense

CS2951K Final Project Proposal

Ning Hou, Lee Painton, Eric Rosen

May 14, 2015

1 Research question

Affect display is the combination of facial, gestural and vocal cues by which persons consciously or unconsciously communicate emotion. Cues such as facial expression, vocal prosody and gestural display are all modes by which the affective state of an individual can be inferred. We are specifically interested in exploring the recognition of confusion in a person by means of their facial landmarks. By confusion we mean the common definition of a mental state where, given a situation, a person either understands it or is confused. We believe that there is a correlation between changes in a person's facial cues and the experience of confusion.

2 Significance

Reliably determining user affect is an open problem in HCI and part of a field called affective computing. The development of affect sensitive intelligent agents would allow computers to interact more effectively with humans in tasks where emotion has an impact, for example learning or driving. Confusion is especially significant during these tasks as it can actively interfere, or even be dangerous in the case of tasks such as operating heavy machinery. Detecting confusion can also be a valuable aid in the diagnoses of medical conditions which may not be immediately obvious to the human observer. Confusion also serves intuitively as a natural perceptual feedback representing the efficacy of an intelligent agent's communication attempts and could be incorporated as part of the reward function in a learning agent. It is our more immediate hope that we

can utilize work in this project to make a Baxter robot aware of confusion in subjects with whom it is interacting.

3 Methodology

In work establishing a framework for machine Emotional Intelligence, Picard et al [?] discuss factors which need to be considered when gathering data for experimental purposes. For our experiment we are interested primarily in event-elicited emotions which arise unconsciously based on the situation. To this end we have designed a quiz of five questions which are intended to provide a spectrum of data. During the course of each question we collect a set of datapoints every 200 milliseconds using an Intel RealSense device which we have attached to a computer and pointed at the quizee. This set includes 13 facial landmarks and 10 emotional features. The facial landmarks are as follows: Left eyebrow raiser, Right eyebrow raiser, Left eyebrow lowerer, Right eyebrow lowerer, Mouth open, Mouth smile, Mouth kiss, Left eye closed, Right eye closed, Eyes turn left, Eyes turn right, Eyes turn up, Eyes turn down and their respective intensities from 0 to 100. The emotional features are a set of 10 statistical features, each comprised of an intensity from 0 to 1 and an evidence rating ranging from -1 to 3. These features include 7 specific emotions and 3 general valence ratings. The 7 emotions are anger, contempt, disgust, fear, joy, sadness, and surprise and the valences are limited to positive, neutral and negative. In addition to these 23 data points we included a single physiological sample in the form of a calculated heart rate.

3.1 Interrater agreement

To establish ground truth we used interrater agreement. Three separate raters viewed videos of the quiz sessions and labelled intervals where they felt the quiz-taker was confused. If two or more raters agreed that a quizee was confused during any given 200ms sample then that frame was considered confusion-positive; otherwise it was confusion-negative.

3.2 Confusion eliciting quiz

To capture the human facial and pulse response to mental state of confusion, we design a short quiz on a scale of questions from easy (less confusing) to hard (more confusing):

1. *What is your name?*
 - Use: Calibrate neutral features.

2. *Who is the President of the United States?*
 - Answer: Barak Obama
 - Use: Easy question measures non-confusing features.
3. *How many fingers am I holding up? (Hold up four)*
 - Answer: Four
 - Use: Easy question measures non-confusing features.
4. *I have two coins totaling 15 cents, one of which is not a nickle. What are the two coins?*
 - Answer: A dime and a nickle
 - Use: Medium question that might seem confusing at first but can be answered after some thinking or clarification. This question measures both confusing and non-confusing features, as well as the transition.
5. *Has anyone really been far enough and decided to use even what they look like?*
 - Answer: Nonsense
 - Use: Intentionally confusing question to measure confusing features.
6. *Theres a dead man in a room surrounded by 53 bicycles. Why is he dead?*
 - Answer: He was caught cheating at cards.
 - Use: Intentionally confusing riddle to measure confusing features.
7. *Make a face of confusion.*
8. *Make a face of understanding.*
 - Use: Extra features of acted features of confusion and non-confusion (understanding).

3.3 Dataset

The dataset consists of two components to the confusion quiz questions in 3.1:

1. The video of the face and upper body of the test subject during the quiz process.
 - Three authors independently annotate the video frames by label of confusion and non-confusion.

- The intersection of annotated confusion frames forms the **baseline** for confusion evaluation and inference.
2. The features detected by Intel RealSense built-in face tracking and emotion modules at each 200 millisecond (ms) frame:
- Pulse [in beats per minute (BPM)]
 - Facial landmarks [on a scale of 0 to 100]:
 - Brow: Raise left [AU1,2], Raise right [AU1,2], Lower left [AU4], Lower right [AU4]
 - Mouth: Smile, Kiss, Open
 - Head: Turn left [AU51], Turn right [AU52], Up [AU53], Down [AU54], Tilt left [AU/M55], Tilt right [AU/M56]
 - Eyes: Turn left [AU/M61], Turn right [AU/M62], Up [AU63], Down [AU64]
 - Emotions [on a scale of 0 to 1 with evidence ratings from -1 to 3]:
 - Primary: Anger, Contempt, Disgust, Fear, Joy, Sadness, Surprise
 - Sentiments: Negative, Positive, Neutral

Some facial landmarks are synonymous with action units as defined by Tian et al [?]. Those have been denoted with the tag [AU]. We could also consider the initial frames of neutral face as AU0. Because mouth features do not overlap with AU definitions, we only use the feature names in our method. However, we use the AU notations here for future use and comparison with other facial action and emotion research. Also, according to EMFACS (Emotional Facial Action Coding System) and FACSaid (Facial Action Coding System Affect Interpretation Dictionary), the seven primary emotions can be detected based on some combinations of action units. It is our current understanding that these methods are used by the Intel RealSense in encoding values for emotional evidence and intensity.

Datasets were collected in the same format under but under two separate scenarios:

1. Conversational: we asked the quiz questions to human subjects.
2. Computerized test: the human subjects took the quiz on computer.

3.4 Naive Bayes

Given our task was simply to label frames either confusion-positive or confusion-negative we felt a simple Naive Bayes classifier to be a worthy attempt. To that end we formulated a 'scaled' bayes classifier which uses partial counts scaled by intensity/evidence ratings for both parameter estimation and classification. We detail this process below.

Consider the features conditionally independent and frames taken at every 200 ms timestamp independent inputs. We formulate the Naive Bayes classifier for confusion mental states: given the feature vector $X = x_1, \dots, x_{28}$ consisting of the 28 features described in Section 3.2.2, we compute $P(\text{confusion}|X)$ by Bayes Rule:

$$P(\text{confusion}|X) = \frac{P(\text{confusion})P(X|\text{confusion})}{P(X)} \quad (1)$$

$$\propto P(\text{confusion})P(X|\text{confusion}) \quad (2)$$

Assuming conditional independence for features in X ,

$$P(\text{confusion}|X) \propto P(\text{confusion}) \prod_{i=1}^{28} P(x_i|\text{confusion}) \quad (3)$$

where

$$P(\text{confusion}) = \frac{\text{intensity scaled count of confusion-positive frames}}{\text{total number of frames} * \text{set of feature intensities}} \quad (4)$$

$$P(x_i|\text{confusion}) = \frac{\text{intensity scaled count of feature } x_i \text{ in confusion class}}{\text{total number of features in confusion class} * \text{feature intensity}} \quad (5)$$

Note that the asterisks above are dot products. Also for the counting of emotional features we only included the feature data from a frame if the evidence rating was greater than 0.

We took questions 1, 5, 7, 8 of Section 3.1 as the training data to compute $P(\text{confusion}|X)$ and evaluate the performance on questions 2, 3, 4, 6 as testing data.

4 Results

4.1 Baseline

We form the baseline of confusion by taking the intersection of three independent sets of annotated frames. The annotation is based on the interrater agreement method established previously. The intersection of this baseline set with the set of all training frames represents the set of all confusion-positive frames in our training data.

We first attempted classification using the full set of all features, then isolated subsets of the features by either eliminating features intuitively or randomly.

4.2 Naive Bayes

Table: Results for scaled Naive Bayes given subsets of data

Unless otherwise stated the threshold for the posterior probability is .50

Feature Subset	Precision	Recall	Accuracy	Fscore
All features	0.208	0.362	0.495	0.265
All features (.70 threshold)	0.218	0.166	0.641	0.188
All features (.80 threshold)	0.251	0.126	0.686	0.168
All features (.90 threshold)	0.300	0.089	0.719	0.138
Without pulse	0.226	0.272	0.583	0.247
Random subset	0.233	0.150	0.663	0.182
Without emotional features	0.198	0.084	0.684	0.118

4.3 Discussion

Our initial attempt at a classifier yielded underwhelming results. We noted that as we increased the threshold

5 Tables and Figures

Figure 1: Sample dataset taken at one 200 ms timestamp

References

- [1] Ashish Kapoor, Selene Mota, and Rosalind W Picard. Towards a learning companion that recognizes affect. In *AAAI Fall symposium*, pages 2–4, 2001.