

Chapter 10

Continuous Analysis of Affect from Voice and Face

Hatice Gunes, Mihalis A. Nicolaou, and Maja Pantic

10.1 Introduction

Human affective behavior is multimodal, continuous and complex. Despite major advances within the affective computing research field, modeling, analyzing, interpreting and responding to human affective behavior still remains a challenge for automated systems as affect and emotions are complex constructs, with fuzzy boundaries and with substantial individual differences in expression and experience [7]. Therefore, affective and behavioral computing researchers have recently invested increased effort in exploring how to best model, analyze and interpret the subtlety, complexity and continuity (represented along a continuum e.g., from -1 to $+1$) of affective behavior in terms of latent dimensions (e.g., arousal, power and valence) and appraisals, rather than in terms of a small number of discrete emotion categories (e.g., happiness and sadness). This chapter aims to (i) give a brief overview of the existing efforts and the major accomplishments in modeling and analysis of emotional expressions in dimensional and continuous space while focusing on open issues and new challenges in the field, and (ii) introduce a representative approach for

H. Gunes (✉) · M.A. Nicolaou · M. Pantic
Imperial College, London, UK
e-mail: h.gunes@imperial.ac.uk

M.A. Nicolaou
e-mail: mihalis@imperial.ac.uk

M. Pantic
e-mail: m.pantic@imperial.ac.uk

H. Gunes
University of Technology Sydney (UTS), Sydney, Australia

M. Pantic
University of Twente, Twente, The Netherlands

47 multimodal continuous analysis of affect from voice and face, and provide exper-
48 imental results using the audiovisual Sensitive Artificial Listener (SAL) Database
49 of natural interactions. The chapter concludes by posing a number of questions that
50 highlight the significant issues in the field, and by extracting potential answers to
51 these questions from the relevant literature.

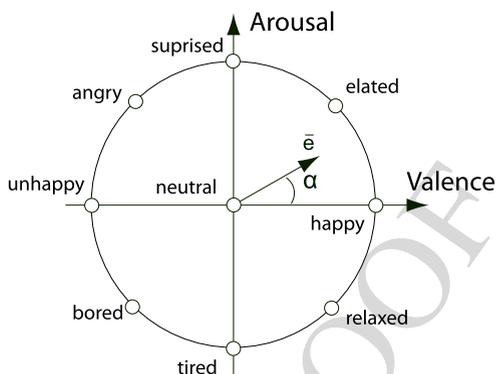
52 The chapter is organized as follows. Section 10.2 describes theories of emotion,
53 Sect. 10.3 provides details on the affect dimensions employed in the literature as
54 well as how emotions are perceived from visual, audio and physiological modal-
55 ities. Section 10.4 summarizes how current technology has been developed, in terms
56 of data acquisition and annotation, and automatic analysis of affect in continuous
57 space by bringing forth a number of issues that need to be taken into account when
58 applying a dimensional approach to emotion recognition, namely, determining the
59 duration of emotions for automatic analysis, modeling the intensity of emotions, de-
60 termining the baseline, dealing with high inter-subject expression variation, defining
61 optimal strategies for fusion of multiple cues and modalities, and identifying appro-
62 priate machine learning techniques and evaluation measures. Section 10.5 presents
63 our representative system that fuses vocal and facial expression cues for dimensional
64 and continuous prediction of emotions in valence and arousal space by employing
65 the bidirectional Long Short-Term Memory neural networks (BLSTM-NN), and in-
66 troduces an output-associative fusion framework that incorporates correlations be-
67 tween the emotion dimensions to further improve continuous affect prediction. Sec-
68 tion 10.6 concludes the chapter.

71 10.2 Affect in Dimensional Space

73 Emotions and affect are researched in various scientific disciplines such as neuro-
74 science, psychology, and cognitive sciences. Development of automatic affect ana-
75 lyzers depends significantly on the progress in the aforementioned sciences. Hence,
76 we start our analysis by exploring the background in emotion theory, perception and
77 recognition.

78 According to research in psychology, three major approaches to affect modeling
79 can be distinguished [31]: categorical, dimensional, and appraisal-based approach.
80 The categorical approach claims that there exist a small number of emotions that
81 are basic, hard-wired in our brain, and recognized universally (e.g. [18]). This the-
82 ory on universality and interpretation of affective nonverbal expressions in terms of
83 basic emotion categories has been the most commonly adopted approach in research
84 on automatic measurement of human affect. However, a number of researchers have
85 shown that in everyday interactions people exhibit non-basic, subtle and rather com-
86 plex affective states like thinking, embarrassment or depression. Such subtle and
87 complex affective states can be expressed via dozens of anatomically possible fa-
88 cial and bodily expressions, audio or physiological signals. Therefore, a single label
89 (or any small number of discrete classes) may not reflect the complexity of the af-
90 fective state conveyed by such rich sources of information [82]. Hence, a number
91 of researchers advocate the use of dimensional description of human affect, where
92

Fig. 10.1 Russel's valence-arousal space. The angle is represented by α while the vector \vec{e} represents the emotion (*point*) as a parameter of arousal and valence. The figure is by courtesy of [77]



93 affective states are not independent from one another; rather, they are related to one
 94 another in a systematic manner (see, e.g., [31, 82, 86]). It is not surprising, therefore,
 95 that automatic affect sensing and recognition researchers have recently started ex-
 96 ploring how to model, analyze and interpret the subtlety, complexity and continuity
 97 (represented along a continuum from -1 to $+1$, without discretization) of affective
 98 behavior in terms of latent dimensions, rather than in terms of a small number of
 99 discrete emotion categories.

100 The most widely used dimensional model is a circular configuration called *Cir-*
 101 *cumplex of Affect* (see Fig. 10.1) introduced by Russell [82]. This model is based on
 102 the hypothesis that each basic emotion represents a bipolar entity being a part of the
 103 same emotional continuum. The proposed poles are arousal (relaxed vs. aroused) and
 104 valence (pleasant vs. unpleasant), as illustrated in Fig. 10.1. Another well-
 105 accepted and commonly used dimensional description is the 3D emotional space of
 106 pleasure—displeasure, arousal—nonarousal and dominance—submissiveness [63],
 107 at times referred to as the *PAD emotion space* [48] or as *emotional primitives* [19].

108 Scherer and colleagues introduced another set of psychological models, referred
 109 to as componential models of emotion, which are based on the appraisal theory
 110 [25, 31, 86]. In the appraisal-based approach emotions are generated through con-
 111 tinuous, recursive subjective evaluation of both our own internal state and the state
 112 of the outside world (relevant concerns/needs) [25, 27, 31, 86]. Despite pioneering
 113 efforts of Scherer and colleagues (e.g., [84]), how to use the appraisal-based ap-
 114 proach for automatic measurement of affect is an open research question as this
 115 approach requires complex, multicomponential and sophisticated measurements of
 116 change. One possibility is to reduce the appraisal models to dimensional models
 117 (e.g., 2D space of arousal-valence).

118 Ortony and colleagues proposed a computationally tractable model of the cog-
 119 nitive basis of emotion elicitation, known as OCC [71]. OCC is now established as
 120 a standard (cognitive appraisal) model for emotions, and has mostly been used in
 121 affect synthesis (in embodied conversational agent design, e.g. [4]).

122 Each approach, categorical or dimensional, has its advantages and disadvantages.
 123 In the categorical approach, where each affective display is classified into a single
 124 category, complex mental states, affective state or blended emotions may be too dif-
 125 ficult to capture.

139 ficult to handle [108]. Instead, in dimensional approach, observers can indicate their
140 impression of each stimulus on several continuous scales. Despite exhibiting such
141 advantages, dimensional approach has received a number of criticisms. Firstly, the
142 usefulness of these approaches has been challenged by discrete emotions theorists,
143 such as Silvan Tomkins, Paul Ekman, and Carroll Izard, who argued that the reduction
144 of emotion space to two or three dimensions is extreme and resulting in loss
145 of information. Secondly, while some basic emotions proposed by Ekman, such as
146 happiness or sadness, seem to fit well in the dimensional space, some basic emotions
147 become indistinguishable (e.g., fear and anger), and some emotions may lie
148 outside the space (e.g., surprise). It also remains unclear how to determine the position
149 of other affect-related states such as confusion. Note, however, that arousal and
150 valence are not claimed to be the only dimensions or to be sufficient to differentiate
151 equally between all emotions. Nonetheless, they have already proven to be useful in
152 several domains (e.g., affective content analysis [107]).

155 10.3 Affect Dimensions and Signals

156
157 An individual's inner emotional state may become apparent by subjective experi-
158 ences (how the person feels), internal/inward expressions (bio signals), and external/
159 outward expressions (audio/visual signals). However, these may be incongruent,
160 depending on the context (e.g., feeling angry and not expressing it outwardly).

161 The contemporary theories of emotion and affect consider appraisal as the most
162 significant component when defining and studying emotional experiences [81], and
163 at the same time acknowledge that emotion is not just appraisal but a complex multi-
164 faceted experience that consists of the following stages (in order of occurrence):

- 166 1. *Cognitive Appraisal*. Only events that have significance for our goals, concerns,
167 values, needs, or well-being elicit emotion.
- 168 2. *Subjective feelings*. The appraisal is accompanied by feelings that are good or
169 bad, pleasant or unpleasant, calm or aroused.
- 170 3. *Physiological arousal*. Emotions are accompanied by autonomic nervous system
171 activity.
- 172 4. *Expressive behaviors*. Emotions are communicated through facial and bodily ex-
173 pressions, postural and voice changes.
- 174 5. *Action tendencies*. Emotions carry behavioral intentions, and the readiness to act
175 in certain ways.

176
177 This multifaceted aspect of affect poses a true challenge to automatic sensing
178 and analysis. Therefore, to be able to deal with these challenges, affect research
179 scientists have ended up making a number of assumptions and simplifications while
180 studying emotions [7, 72]. These assumptions can be listed.

- 181 1. *Emotions are on or off at any particular point in time*. This assumption has im-
182 plications on most data annotation procedures where raters label a user's ex-
183 pressed emotion as one of the basic emotion categories or a specific point in a
184

185 dimensional space. The main issue with this assumption is that the boundaries
186 for defining the expressed emotion as on or off are usually not clear.

- 187 2. *Emotion is a state that the subject does not try to actively change or alleviate.*
188 This is a common assumption during the data acquisition process where the sub-
189 jects are assumed to have a simple response to the provided stimulus (e.g., while
190 watching a clip or interacting with an interface). However, such simple passive
191 responses do not usually hold during daily human–computer interactions. Peo-
192 ple generally regulate their affective states caused by various interactions (e.g.,
193 an office user logging into Facebook to alleviate his boredom).
- 194 3. *Emotion is not affected by situation or context.* This assumption pertains to most
195 of the past research work on automatic affect recognition where emotions have
196 been mostly investigated in laboratory settings, outside of a social context. How-
197 ever, some emotional expressions are displayed only during certain context (e.g.,
198 pain).

199 Affect research scientists have made the following simplifications while studying
200 emotions [7, 72]:

- 201 1. *Emotions do occur in asynchronous communication* (e.g., via a prerecorded
202 video/sound from a sender to a receiver). This simplification does not hold in re-
203 ality as human nonverbal expressive communication occurs mostly face-to-face.
- 204 2. *Interpersonal emotions do arise from communications with strangers* (e.g., lab-
205 oratory studies where people end up communicating with people they do not
206 know). This simplification is unrealistic as people tend to be less expressive with
207 people they do not know on an interpersonal level. Therefore, an automatic sys-
208 tem designed using such communicative settings is expected to be much less
209 sensitive to its user’s realistic expressions.
210

211 Overall, these assumptions and simplifications are far from reality. However, they
212 have paved the initial but crucial way for automatic affect recognizers that attempt to
213 analyze both the felt (e.g., [9, 10, 59]) and the internally or the externally expressed
214 (e.g., [50, 54]) emotions.

217 **10.3.1 Affect Dimensions**

218
219 Despite the existence of various emotion models described in Sect. 10.2, in auto-
220 matic measurement of dimensional and continuous affect, valence (how positive or
221 negative the affect is), activation (how excited or apathetic the affect is), power (the
222 sense of control over the affect), and expectation (the degree of anticipating or being
223 taken unaware) appear to make up the four most important affect dimensions [25].
224 Although ideally the intensity dimension could be derived from the other dimen-
225 sions, to guarantee a complete description of affective coloring, some researchers
226 include intensity (how far a person is away from a state of pure, cool rationality)
227 as the fifth dimension (e.g., [62]). Solidarity, antagonism and agreement have also
228 been in the list of dimensions investigated [13]. Overall, search for optimal low-
229 dimensional representation of affect remains open [25].
230

10.3.2 Visual Signals

Facial actions (e.g., pulling eyebrows up) and facial expressions (e.g., producing a smile), and to a much lesser extent bodily postures (e.g., head bent backwards and arms raised forwards and upwards) and expressions (e.g., head nod), form the widely known and used visual signals for automatic affect measurement. Dimensional models are considered important in this task as a single label may not reflect the complexity of the affective state conveyed by a facial expression, body posture or gesture. Ekman and Friesen [17] considered expressing discrete emotion categories via face, and communicating dimensions of affect via body as more plausible.

A number of researchers have investigated how to map various visual signals onto emotion dimensions. For instance, Russell [82] mapped the facial expressions to various positions on the two-dimensional plane of arousal-valence, while Cowie et al. [13] investigated the emotional and communicative significance of head nods and shakes in terms of arousal and valence dimensions, together with dimensional representation of solidarity, antagonism and agreement.

Although in a stricter sense not seen as part of the visual modality, motion capture systems have also been utilized for recording the relationship between body posture and affect dimensions (e.g., [57, 58]). For instance, Kleinsmith et al. [58] identified that scaling, arousal, valence, and action tendency were the affective dimensions used by human observers when discriminating between postures. They also reported that low-level posture features such as orientation (e.g., orientation of shoulder axis) and distance (e.g., distance between left elbow and right shoulder) appear to help in effectively discriminating between the affective dimensions [57, 58].

10.3.3 Audio Signals

Audio signals convey affective information through explicit (linguistic) messages, and implicit (acoustic and prosodic) messages that reflect the way the words are spoken. There exist a number of works focusing on how to map audio expression to dimensional models. Cowie et al. used valence-activation space (similar to valence-arousal) to model and assess affect from speech [11, 12]. Scherer and colleagues have also proposed how to judge emotional effects on vocal expression, using the appraisal-based theory [31].

In terms of affect recognition from audio signals the most reliable finding is that pitch appears to be an index into arousal [7]. Another well-accepted finding is that mean of the fundamental frequency (F0), mean intensity, speech rate, as well as pitch range [46], “blaring” timbre [14] and high-frequency energy [85] are positively correlated with the arousal dimension. Shorter pauses and inter-breath stretches are indicative of higher activation [99].

There is relatively less evidence on the relationship between certain acoustic parameters and other affect dimensions such as valence and power. Vowel duration and power dimension in general, and lower F0 and high power in particular, appear

277 to have correlations. Positive valence seems to correspond to a faster speaking rate,
278 less high-frequency energy, low pitch and large pitch range [85] and longer vowel
279 durations. A detailed literature summary on these can be found in [87] and [88].
280

281 282 **10.3.4 Bio Signals** 283

284 The bio signals used for automatic measurement of affect are galvanic skin response
285 that increases linearly with a person's level of arousal [9], electromyography (fre-
286 quency of muscle tension) that is correlated with negatively valenced emotions [41],
287 heart rate that increases with negatively valenced emotions such as fear, heart rate
288 variability that indicates a state of relaxation or mental stress, and respiration rate
289 (how deep and fast the breath is) that becomes irregular with more aroused emotions
290 like anger or fear [9, 41].
291

292 Measurements recorded over various parts of the brain including the amygdala
293 also enable observation of the emotions felt [79]. For instance, approach or with-
294 drawal response to a stimulus is known to be linked to the activation of the left or
295 right frontal cortex, respectively.

296 A number of studies also suggest that there exists a correlation between in-
297 creased blood perfusion in the orbital muscles and stress levels for human beings.
298 This periorbital perfusion can be quantified through the processing of thermal video
299 (e.g., [102]).
300

301 302 **10.4 Overview of the Current Technology** 303

304 This section provides a brief summary of the current technology by describing how
305 affective data are acquired and annotated, and how affect analysis in continuous
306 space is achieved.
307

308 309 **10.4.1 Data Acquisition and Annotation** 310

311 Cameras are used for acquisition of face and bodily expressions, microphones are
312 used for recording audio signals, and thermal (infrared) cameras are used for record-
313 ing blood flow and changes in skin temperature. 3D affective body postures or
314 gestures can alternatively be recorded by utilizing motion capture systems (e.g.,
315 [57, 58]). In such scenarios, the actor is dressed in a suit with a number of markers
316 on the joints and body segments, while each gesture is captured by a number of cam-
317 eras and represented by consecutive frames describing the position of the markers
318 in the 3D space. This is illustrated in Fig. 10.2 (second and third rows).
319

320 In the bio signal research context, the subject being recorded usually wears a
321 headband or a cap on which electrodes are mounted, a clip sensor, or touch type
322

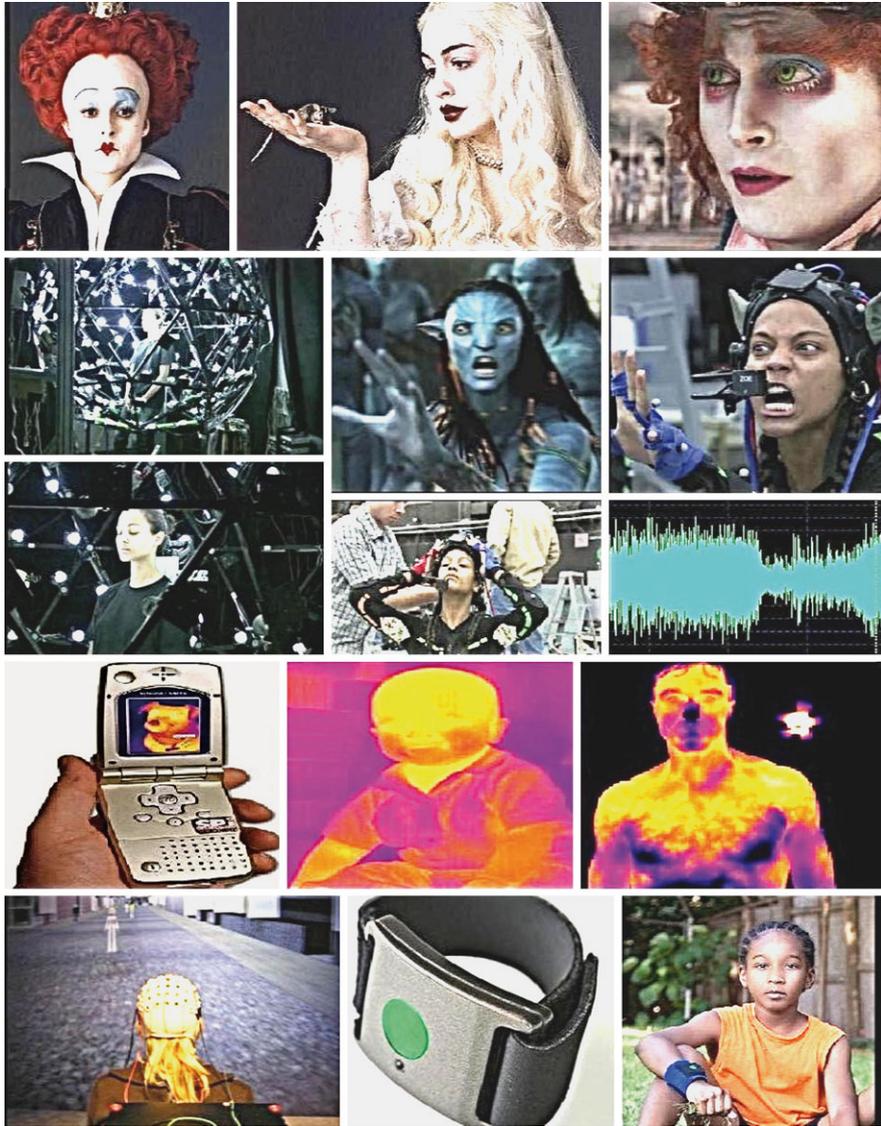


Fig. 10.2 Examples of sensors used in multimodal affective data acquisition: (*1st row*) camera for visible imagery (face and body), (*2nd & 3rd rows*) facial and body motion capture, and audio signals (used for animation and rendering), (*4th row*) infrared camera for thermal imagery, and (*5th row*) various means for recording bio signals (brain signals, heart and respiration rate, etc.)

electrodes (see Fig. 10.2, last row). The subject is then stimulated with emotionally-evocative images or sounds. Acquiring affect data without subjects' knowledge is strongly discouraged and the current trend is to record spontaneous data in more

369 constrained conditions such as an interview (e.g., [10]) or interaction (e.g., [62])
370 setting, where subjects are still aware of placement of the sensors and their locations.

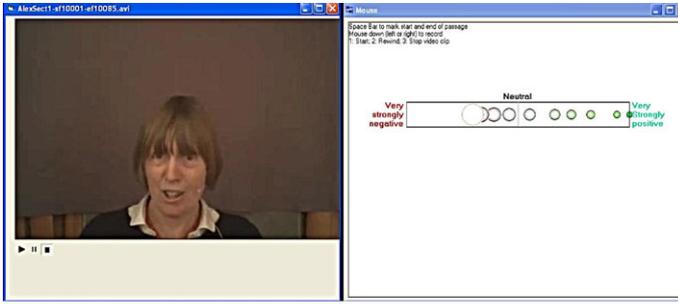
371 Annotation of the affect data is usually done separately for each modality, as-
372 suming independency between the modalities. A major challenge is the fact that
373 there is no coding scheme that is agreed upon and used by all researchers in the field
374 that can accommodate all possible communicative cues and modalities. In general,
375 the Feeltrace annotation tool is used for annotating the external expressions (audio
376 and visual signals) with continuous traces (impressions) in the dimensional space.
377 Feeltrace allows coders to watch the audiovisual recordings and move their cursor,
378 within the 2-dimensional emotion space (valence and arousal) confined to $[-1, +1]$,
379 to rate their impression about the emotional state of the subject [11] (see the illustra-
380 tion in Fig. 10.3(a)). For annotating the internal expressions (bio signals), the level
381 of valence and arousal is usually extracted from subjective experiences (subjects'
382 own responses) (e.g., [59, 79]) due to the fact that feelings, induced by an image
383 or sound, can be very different from subject to subject. The Self Assessment Man-
384 nequin (SAM) [60], illustrated in Fig. 10.3(b), is the most widely used means for
385 self assessment.

386 When discretized dimensional annotation is adopted (as opposed to continuous
387 one), researchers seem to use different intensity levels: either a ten-point Likert scale
388 (e.g., 0-low arousal, 9-high arousal) or a range between -1.0 and 1.0 (divided into
389 a number of levels) [37]. The final annotation is usually calculated as the mean of
390 the observers' ratings. However, whether this is the best way of obtaining ground-
391 truth labels of emotional data is still being discussed. Overall, individual coders
392 may vary in their appraisal of what is happening in the scene, in their judgment of
393 the emotional behavior of the target individual, in their understanding of the terms
394 'positive emotion' and 'negative emotion' and in their movement of the computer
395 mouse to translate their rating into a point on the onscreen scale. Furthermore, recent
396 findings in dynamic emotional behavior coding indicate that the temporal pattern of
397 ratings appears similar across cultures but that there exist significant differences in
398 the intensity levels at which participants from different cultural backgrounds rate
399 the emotional behaviors [96]. Therefore, how to obtain and use rich emotional data
400 annotations, from multiple and multi-cultural raters, needs serious consideration.

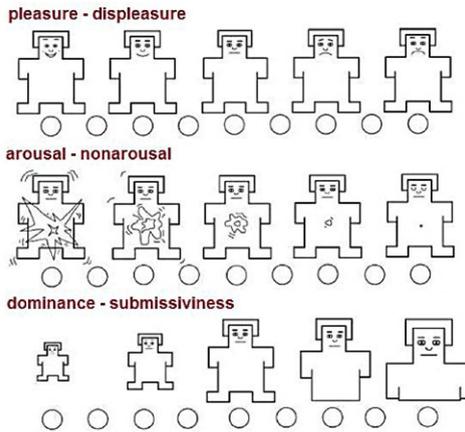
401 402 403 ***10.4.2 Automatic Dimensional Affect Prediction and Recognition*** 404

405 After affect data have been acquired and annotated, representative and relevant fea-
406 tures need to be extracted prior to the automatic measurement of affect in dimen-
407 sional and continuous space. The feature extraction techniques used for each com-
408 municative source are similar to the previous works (reviewed in [40]) adopting a
409 categorical approach to affect recognition.

410 In dimensional affect analysis emotions are represented along a continuum. Con-
411 sidering this, systems that target automatic dimensional affect measurement should
412 be able to predict the emotions continuously. However, most of the automatic recog-
413 nition systems tend to simplify the problem by quantizing the continuous labels into
414



(a)



(b)

Fig. 10.3 Illustration of (a) the Feeltrace annotation tool [11], and (b) the Self Assessment Mannequin (SAM) [60]

a finite number of discrete levels. Hence, the most commonly employed strategy in automatic dimensional affect prediction is to reduce the continuous prediction problem to a two-class recognition problem (positive vs. negative or active vs. passive classification; e.g., [66, 92]) or a four-class recognition problem (classification into the quadrants of 2D V-A space; e.g., [8, 26, 29, 47, 106]).

For example, Kleinsmith and Bianchi-Berthouze discriminate between high–low, high–neutral and low–neutral affective dimensions [57], while Wöllmer et al. quantize the V-A dimensions of the SAL database into either 4 or 7 levels, and then use Conditional Random Fields (CRFs) to predict the quantized labels [105]. Attempts for discriminating between more coarse categories, such as positive vs. negative [66], and active vs. passive [8] have also been attempted. Of these, Caridakis et al. [8] uses the SAL database, combining auditive and visual modalities. Nicolaou et al. focus on audiovisual classification of spontaneous affect into negative or positive emotion categories using facial expression, shoulder and audio cues, and utilizing 2- and 3-chain coupled Hidden Markov Models and likelihood space classification to

461 fuse multiple cues and modalities [66]. Kanluan et al. combine audio and visual cues
462 for affect recognition in V-A space by fusing facial expression and audio cues, using
463 Support Vector Machines for Regression (SVR) and late fusion with a weighted linear
464 combination [50]. The labels used have been discretized on a 5-point scale in the
465 range of $[-1, +1]$ for each emotion dimension. The work presented in [106] utilizes
466 a hierarchical dynamic Bayesian network combined with BLSTM-NN performing
467 regression and quantizing the results into four quadrants (after training).

468 As far as actual continuous dimensional affect prediction (without quantization)
469 is concerned, there exist a number of methods that deal exclusively with speech
470 (i.e., [33, 105, 106]). The work by Wöllmer et al. uses the SAL Database and Long
471 Short-Term Memory neural networks and Support Vector Machines for Regression
472 (SVR) [105]. Grimm and Kroschel use the Vera am Mittag database [35] and SVRs,
473 and compare their performance to that of the distance-based fuzzy k-Nearest Neighbor
474 and rule-based fuzzy-logic estimators [33]. The work by Espinosa et al. also use
475 the Vera am Mittag database [35] and examine the importance of different groups
476 of speech acoustic features in the estimation of continuous PAD dimensions [19].

477 Currently, there are also a number of works focusing on dimensional and continuous
478 prediction of emotions from the visual modality [39, 56, 69]. The work by
479 Gunes and Pantic focuses on dimensional prediction of emotions from spontaneous
480 conversational head gestures by mapping the amount and direction of head motion,
481 and occurrences of head nods and shakes into arousal, expectation, intensity, power
482 and valence level of the observed subject using SVRs [39]. Kipp and Martin in [56]
483 investigated (without performing automatic prediction) how basic gestural form features
484 (e.g., preference for using left/right hand, hand shape, palm orientation, etc.)
485 are related to the single PAD dimensions of emotion. The work by Nicolaou et al.
486 focuses on dimensional and continuous prediction of emotions from naturalistic facial
487 expressions within an Output-Associative Relevance Vector Machine (RVM)
488 regression framework by learning non-linear input and output dependencies inherent
489 in the affective data [69].

490 More recent works focus on dimensional and continuous prediction of emotions
491 from multiple modalities. For instance, Eyben et al. [21] propose a string-based
492 approach for fusing the behavioral events from visual and auditive modalities (i.e.,
493 facial action units, head nods and shakes, and verbal and nonverbal audio cues) to
494 predict human affect in a continuous dimensional space (in terms of arousal, expectation,
495 intensity, power and valence dimensions). Although automatic affect analyzers based
496 on physiology end up using multiple signal sources, explicit fusion of multimodal
497 data for continuous modeling of affect utilizing dimensional models of emotion is
498 still relatively unexplored. For instance, Khalili and Moradi propose multimodal
499 fusion of brain and peripheral signals for automatic recognition of three emotion
500 categories (positively excited, negatively excited and calm) [52]. Their results
501 show that, for the task at hand, EEG signals seem to perform better than other
502 physiological signals, and nonlinear features lead to better understanding of the felt
503 emotions. Another representative approach is that of Gilroy et al. [28] that propose
504 a dimensional multimodal fusion scheme based on the power-arousal-PAD space
505 to support detection and integration of spontaneous affective behavior of users (in
506

507 terms of audio, video and attention events) experiencing arts and entertainment. Un-
508 like many other multimodal approaches (e.g., [8, 50, 66]), the ground truth in this
509 work is obtained by measuring Galvanic Skin Response (GSR) as an independent
510 measure of arousal.

511 For further details on the aforementioned systems, as well as on systems that
512 deal with dimensional affect recognition from a single modality or cue, the reader
513 is referred to [37, 38, 109].

514 515 516 **10.4.3 Challenges and Prospects** 517

518
519 The summary provided in the previous section reflects that automatic dimensional
520 affect recognition is still in its pioneering stage [34, 37, 38, 91, 105]. There are a
521 number of challenges which need to be taken into account when applying a dimen-
522 sional approach to affect prediction and advancing the current state of the art.

523 *The interpretation accuracy* of expressions and physiological responses in terms
524 of continuous emotions is very challenging. While visual signals appear to be bet-
525 ter for interpreting valence, audio signals seem to be better for interpreting arousal
526 [33, 68, 100, 105]. A thorough comparison between all modalities would indeed
527 provide a better understanding of which emotion dimensions are better predicted
528 from which modalities (or cues).

529 *Achieving inter-observer agreement* is one of the most challenging issues in
530 dimension-based affect modeling and analysis. To date, researchers have mostly
531 chosen to use self-assessments (subjective experiences, e.g. [41]) or the mean
532 (within a predefined range of values) of the observers' ratings (e.g. [57]). Although
533 it is difficult to self-assess arousal, it has been reported that using classes gener-
534 ated from self-assessment of emotions facilitate greater accuracy in recognition
535 (e.g., [9]). This finding results from a study on automatic analysis of physiologi-
536 cal signals in terms of A-V emotion space. It remains unclear whether the same
537 holds independently of the utilized modalities and cues. Modeling inter-observer
538 agreement levels within automatic affect analyzers and finding which signals bet-
539 ter correlate with self assessment and which ones better correlate with independent
540 observer assessment remain unexplored.

541 *The window size* to be used to achieve optimal affect prediction is another is-
542 sue that the existing literature does not provide a unique answer to. Current affect
543 analyzers employ various window sizes depending on the modality, e.g., 2–6 sec-
544 onds for speech, 3–15 seconds for bio signals [54]. For instance, when measuring
545 affect from heart rate signals, analysis should not be done on epochs of less than
546 a minute [6]. A time window of 50 s appears to be also necessary to accurately
547 monitor mental stress in realistic settings [83]. There is no consensus on how the
548 efficiency of such a choice should be evaluated. On one hand achieving real-time
549 affect prediction requires a small window size to be used for analysis (i.e., a few
550 seconds, e.g. [10]), while on the other hand obtaining a reliable prediction accuracy
551 requires long(er)-term monitoring [6, 83]. For instance, Chanel et al. [10] conducted
552

553 short-term analysis of emotions (i.e., time segments of 8 s) in valence and arousal
554 space using EEG and peripheral signals in a self-induction paradigm. They reported
555 large differences in accuracy between the EEG and peripheral features which may
556 be due to the fact that the 8 s length of trials may be too short for a complete activa-
557 tion of peripheral signals while it may be sufficient for EEG signals.

558 *Measuring the intensity* of expressed emotion appears to be modality dependent.
559 The way the intensity of an emotion is apparent from physiological data may be
560 different from the way it is apparent from visual data. Moreover, little attention has
561 been paid so far to whether there are definite boundaries along the affect continuum
562 to distinguish between various levels or intensities. Currently intensity is measured
563 by quantizing the affect dimensions into arbitrary number of levels such as neutral,
564 low and high (e.g., [57, 59, 105]). Separate models are then built to discriminate
565 between pairs of affective dimension levels, for instance, low vs. high, low vs. neu-
566 tral, etc. Generalizing intensity analysis across different subjects is a challenge yet
567 to be researched as different subjects express different levels of emotions in the
568 same situation. Moreover, recent research findings indicate that there also exist sig-
569 nificant differences in the intensity levels at which coders from different cultural
570 backgrounds rate emotional behaviors [96].

571 *The Baseline problem* is another major challenge in the field. For physiologi-
572 cal signals (bio signals) this refers to the problem of finding a condition against
573 which changes in measured physiological signals can be compared (a state of calm-
574 ness) [65]. For the audio modality this is usually achieved by segmenting the record-
575 ings into turns using energy based voice activity detection and processing each turn
576 separately (e.g., [105]). For visual modality the aim is to find a frame in which the
577 subject is expressionless and against which changes in subject's motion, pose, and
578 appearance can be compared. This is achieved by manually segmenting the record-
579 ings, or by constraining the recordings to have the first frame containing a neutral
580 expression (see, e.g., [66, 67, 75]). Yet, as pointed out by Levenson in [61], emo-
581 tion is rarely superimposed upon a prior state of *rest*; instead, emotion occurs most
582 typically when the organism is in some prior activation. Hence, enforcing existence
583 of expressionless state in each recording or manually segmenting recordings so that
584 each segment contains a baseline expression are strong, unrealistic constrains. This
585 remains a great challenge in automatic analysis, which typically relies on existence
586 of a baseline for analysis and processing of affective information.

587 *Generalization capability* of automatic affect analyzers across subjects is still a
588 challenge in the field. Kulic and Croft [59] reported that for bio signal based af-
589 fect measurement, subjects seem to vary not only in terms of response amplitude
590 and duration, but for some modalities, a number of subjects show no response at all.
591 This makes generalization over unseen subjects a very difficult problem. A common
592 way of measuring affect from bio signals is doing it for each participant separately
593 (without computing baseline), e.g. [10]. When it comes to other modalities, most of
594 the works in the field report mainly on subject-dependent dimensional affect mea-
595 surement and recognition due to limited number of subjects and limited amount of
596 data (e.g., [39, 68, 69, 105]).

597 *Modality fusion* refers to combining and integrating all incoming unimodal
598 events into a single representation of the affect expressed by the user. When it

599 comes to integrating multiple modalities, the major issues are: (i) when to integrate
600 the modalities (at what abstraction level to do the fusion), (ii) how to integrate the
601 modalities (which criteria to use), (iii) how to deal with the increased number of
602 features due to fusion, (iv) how to deal with the asynchrony between the modalities
603 (e.g., if video is recorded at 25 Hz, audio is recorded at 48 kHz while EEG is
604 recorded at 256–512 Hz), and (v) how to proceed with fusion when there is con-
605 flicting information conveyed by the modalities. Typically, multimodal data fusion
606 is either done at the feature level (in a maximum likelihood estimation manner) or
607 at the decision level (when most of the joint statistical properties may have been
608 lost). Feature-level fusion is obtained by concatenating all the features from multiple
609 cues into one feature vector which is then fed into a machine learning technique.
610 In the decision-level data fusion, the input coming from each modality/cue is modeled
611 independently, and these single-cue and single-modality based recognition results
612 are combined in the end. Since humans display multi-cue and multimodal
613 expressions in a complementary and redundant manner, the assumption of conditional
614 independence between modalities and cues in decision-level fusion can result in loss
615 of information (i.e. mutual correlation between the modalities). Therefore, model-
616 level fusion has been proposed as an alternative approach for fusing multimodal
617 affect data (e.g., [75]). Despite such efforts in the discrete affect recognition field
618 (reviewed in [40, 109]), these issues remain yet to be explored for dimensional
619 and continuous affect prediction.

620 *Machine learning techniques* used for dimensional and continuous affect measure-
621 ment should be able to produce continuous values for the target dimensions. Overall,
622 there is no agreement on how to model dimensional affect space (continuous vs. quantized)
623 and which machine learning technique is better suited for automatic, multimodal, continuous
624 affect analysis using a dimensional representation. Recognition of quantized dimensional
625 labels is obtained via classification while continuous prediction is achieved by regression.
626 Conditional Random Fields (CRF) and Support Vector Machines (SVM) have mostly been
627 used for quantized dimensional affect recognition tasks (e.g., [105]). Some of the schemes
628 that have been explored for the task of prediction are Support Vector Machines for Regression
629 (SVR) (e.g., [39]) and Long Short-Term Memory Recurrent Networks (LSTM-RNN).
630 The design of emotion-specific classification schemes that can handle multimodal and
631 spontaneous data is one of the most important issues in the field. In accordance with
632 this, Kim and Andre propose a novel scheme of emotion-specific multilevel dichotomous
633 classification (EMDC) using the property of the dichotomous categorization in the 2D
634 emotion model and the fact that arousal classification yields a higher correct classification
635 ratio than valence classification (or direct multiclass classification) [55]. They apply
636 this scheme on classification of four emotions (positive/high arousal, negative/high
637 arousal, negative/low arousal and positive/low arousal) from physiological signals
638 recorded while subjects were listening to music. How to create such emotion-specific
639 schemes for dimensional and continuous prediction of emotions from other modalities
640 and cues should be investigated further.

641 *Evaluation measures* applicable to categorical affect recognition are not directly
642 applicable to dimensional approaches. Using the Mean Squared Error (MSE) between
643 the predicted and the actual values of arousal and valence, instead of the
644

645 recognition rate (i.e., percentage of correctly classified instances) is the most commonly
646 used measure by related work in the literature (e.g., [50, 105]). However,
647 using MSE might not be the best way to evaluate the performance of dimensional
648 approaches to automatic affect measurement and prediction. Therefore, the correlation
649 coefficient that evaluates whether the model has managed to capture patterns
650 inhibited in the data at hand is also employed by several studies (e.g., [50, 67])
651 together with MSE. Overall, however, how to obtain optimal evaluation metrics for
652 continuous and dimensional emotion prediction remains an open research issue [37].
653 Generally speaking, the performance of an automatic analyzer can be modeled and
654 evaluated in an *intrinsic* and an *extrinsic* manner (as proposed for face recognition
655 in [103]). The intrinsic performance and its evaluation depend on the intrinsic components
656 such as the dataset chosen for the experiments and the machine learning algorithms
657 (and their parameters) utilized for prediction. The extrinsic performance and evaluation
658 instead depend on the extrinsic factors such as (temporal/spatial) resolution of the
659 multimodal data and recording conditions (e.g., illumination, occlusions, noise, etc.).
660 Future research in continuous affect prediction should analyze the relevance and prospects
661 of the aforementioned performance components, and how they could be applied to
662 continuous prediction of affect.
663

664 665 **10.4.4 Applications** 666

667
668 Various applications have been using the dimensional (both quantized and continuous)
669 representation and prediction of emotions, ranging from human–computer
670 (e.g., Sensitive Talking Heads [45], Sensitive Artificial Listeners [89, 90], spatial
671 attention analysis [95], arts installations [104]) and human–robot interaction (e.g.,
672 humanoid robotics [5, 51]), clinical and biomedical studies (e.g., stress/pain monitoring
673 [36, 64, 101], autism-related assistive technology), learning and driving environments
674 (e.g., episodic learning [22], affect analysis in the car [20]), multimedia
675 (e.g., video content representation and retrieval [53, 98] and personalized affective
676 video retrieval [97]), and entertainment technology (e.g., gaming [80]). These indicate
677 that affective computing has matured enough to have a presence and measurable
678 impact in our lives. There are also spin off companies emerging out of collaborative
679 research at well-known universities (e.g., Affectiva [1] established by R. Picard and
680 colleagues of MIT Media Lab).
681

682 683 **10.5 A Representative System: Continuous Analysis of Affect 684 from Voice and Face** 685

686
687 The review provided in the previous sections indicates that currently there is a shift
688 toward subtle, continuous, and context-specific interpretations of affective displays
689 recorded in naturalistic settings, and toward multimodal analysis and recognition of
690

691 human affect. Converging with this shift, in this section we present a representative
692 approach that: (i) fuses facial expression and audio cues for dimensional and contin-
693 uous prediction of emotions in valence and arousal space, (ii) employs the bidirec-
694 tional Long Short-Term Memory neural networks (BLSTM-NNs) for the prediction
695 task, and (iii) introduces an output-associative fusion framework that incorporates
696 correlations between the emotion dimensions to further improve continuous predic-
697 tion of affect.

698 The section starts with the description of the naturalistic database used in the
699 experimental studies. Next, data pre-processing, audio and facial feature extraction
700 and tracking procedures, as well as the affect prediction process are explained.

701 702 703 **10.5.1 Dataset**

704
705 We use the Sensitive Artificial Listener Database (SAL-DB) [16] that contains spon-
706 taneous data collected with the aim of capturing the audiovisual interaction between
707 a human and an operator undertaking the role of a SAL character (e.g., an avatar).
708 The SAL characters intend to engage the user in a conversation by paying atten-
709 tion to the user's emotions and nonverbal expressions. Each character has its own
710 emotionally defined personality: Poppy is happy, Obadiah is gloomy, Spike is angry,
711 and Prudence is pragmatic. During an interaction, the characters attempt to create an
712 emotional workout for the user by drawing her/him toward their dominant emotion,
713 through a combination of verbal and nonverbal expressions.

714 The SAL database contains audiovisual sequences recorded at a video rate of
715 25 fps (352×288 pixels) and at an audio rate of 16 kHz. The recordings were
716 made in a lab setting, using one camera, a uniform background and constant light-
717 ing conditions. The SAL data have been annotated manually. Although there are
718 approximately 10 hours of footage available in the SAL database, V-A annotations
719 have only been obtained for two female and two male subjects. We used this portion
720 for our experiments.

721 722 723 **10.5.2 Data Pre-processing and Segmentation**

724
725 The data pre-processing and segmentation stage consists of (i) determining ground
726 truth by maximizing inter-coder agreement, (ii) detecting frames that capture the
727 transition *to* and *from* an emotional state, and (iii) automatic segmentation of spon-
728 taneous audiovisual data. We provide a brief summary of these in the following
729 sections. For a detailed description of these procedures the reader is referred to [67].

730 731 732 **10.5.2.1 Annotation Pre-processing**

733
734 The SAL data have been annotated by a set of coders who provided continuous
735 annotations with respect to valence and arousal dimensions using the Feeltrace an-
736

737 notation tool [11], as explained in Sect. 10.4.1. Feeltrace allows coders to watch
 738 the audiovisual recordings and move their cursor, within the 2-dimensional emotion
 739 space (valence and arousal) confined to $[-1, +1]$, to rate their impression about the
 740 emotional state of the subject.

741 Annotation pre-processing involves dealing with the issue of missing values
 742 (interpolation), grouping the annotations that correspond to one video frame to-
 743 gether (binning), determining normalization procedures (normalization) and extract-
 744 ing statistics from the data in order to obtain segments with a baseline and high
 745 inter-coder agreement (statistics and metrics).

747 **Interpolation** In order to deal with the issue of missing values, similar to other
 748 works reporting on data annotated in continuous dimensional spaces (e.g., [105]),
 749 we interpolated the actual annotations at hand. We used piecewise cubic interpola-
 750 tion as it preserves the monotonicity and the shape of the data.

752 **Binning** Binning refers to grouping and storing the annotations together. As a
 753 first step the measurements of each coder c are binned separately. Since we aim at
 754 segmenting video files, we generate bins which are equivalent to one video frame f .
 755 This is equivalent to a bin of 0.04 seconds (SAL-DB was recorded at a rate of
 756 25 frames/s). The fields with no annotation are assigned a ‘not a number’ (NaN)
 757 identifier.

760 **Normalization** The A-V measurements for each coder are not in total agreement,
 761 mostly due to the variance in human coders’ perception and interpretation of emo-
 762 tional expressions. Thus, in order to deem the annotations comparable, we need to
 763 normalize the data. We experimented with various normalization techniques. After
 764 extracting the videos and inspecting the superimposed ground-truth plots, we opted
 765 for local normalization (normalizing each coder file for each session). This helps us
 766 avoid propagating noise in cases where one of the coders is in large disagreement
 767 with the rest (where a coder has a very low correlation with respect to the rest of
 768 the coders). Locally normalizing to zero mean produces the smallest mean squared
 769 error (MSE) both for valence (0.046) and arousal (0.0551) dimensions.

771 **Statistics and Metrics** We extract two useful statistics from the annotations: cor-
 772 relation and agreement. We start the analysis by constructing vectors of pairs of
 773 coders that correspond to each video session, e.g., when we have a video session
 774 where four coders have provided annotations, this gives rise to six pairs. For each
 775 of these pairs we extract the correlation coefficient between the valence (val) values
 776 of each pair, as well as the level of agreement in emotion classification in terms of
 777 positive or negative. We define the agreement metric by
 778

$$780 \quad AGR = \frac{\sum_{f=0}^n e(c_i(f).val, c_j(f).val)}{|frames|}, \quad (10.1)$$

where $c_i(f).val$ stands for the valence value annotated by coder c_i at frame f . Function e is defined as

$$e(i, j) = \begin{cases} 1 & \text{if } (sign(i) = sign(j)), \\ 0 & \text{else.} \end{cases}$$

In these calculations we do not consider the NaN values to avoid negatively affecting the results. After these metrics are calculated for each pair, each coder is assigned the average of the results of all pairs that the coder has participated in. We choose the Pearson's Correlation (COR) as the metric to be used in the automatic segmentation process as it appears to be stricter than agreement (AGR) providing better comparison amongst the coders.

10.5.2.2 Automatic Segmentation

The segmentation stage consists of producing negative and positive audiovisual segments with a temporal window that contains an offset before and after (i.e., the baseline) the displayed expression. For instance, for capturing negative emotional states, if we assume that the transition *from* non-negative *to* negative emotional state occurs at time t (in seconds), we would have a window of $[t - 1, t, t', t' + 1]$ where t' seconds is when the emotional state of the subject turns to non-negative again. The procedure is completely analogous for positive emotional states.

Detecting and Matching Crossovers For an input coder c , the crossing over from one emotional state to the other is detected by examining the valence values and identifying the points where the sign changes. Here a modified version of the sign function is used, it returns 1 for values that are higher than 0 (a value of 0 valence is never encountered in the annotations), -1 for values that are less than zero, and 0 for NaN values. We accumulate all crossover points for each coder, and return the set of crossovers *to-a-positive* and *to-a-negative* emotional state. The set of crossovers is then used for matching crossovers across coders. For instance, if a session has annotations from four coders, the frame (f) where each coder detects the crossover is not the same for all coders (for the session in question). Thus, we have to allow an offset for the matching process. This procedure searches the crossovers detected by the coders and then accepts the matches where there is less than the predefined offset (time) difference between the detections. When a match is found, we remove the matched crossovers and continue with the rest. The existence of different combinations of crossovers which may match using the predefined offset poses an issue. By examining the available datasets, we decided to maximize the number of coders participating in a matched crossover set rather than minimizing the temporal distances between the participating coders. The motivations for this decision are as follows: (i) if more coders agree on the crossover, the reliability of the ground truth produced will be higher, and (ii) the offset amongst the resulting matches is on average quite small (<0.5 s) when considering only the number of participating coders. We disregard cases where only one coder detects a crossover due to lack of agreement between coders.

Segmentation Driven by Matched Crossovers In order to illustrate how the crossover frame decision (for each member of the set) is made, let us assume that for *to-a-negative* transition a coder detects a crossover at frame 2, while the other coder detects a crossover at frame 4. If the frames are averaged to the nearest integer, then we can assume that the crossover happens at frame 3. In this case we have only 2 coders agreeing, we use the *correlation* metric in order to weight their decision and determine the crossover point. This provides a measurement of the relative importance of the annotations for each coder and propagates information from the other two coders not participating in the match. In order to capture 0.5 s before the transition window, the number of frames corresponding to the predefined offset are subtracted from the *start frame*. The ground-truth values for valence are retrieved by incrementing the initial frame number where each crossover was detected by the coders. Again, following the previous example, this means that we consider frame 2 of coder 1 and frame 4 of coder 2 to provide ground-truth values for frame 3 (the average of 2 and 4). This gives us an averaged valence value. Then, the frame 4 valence value (ground truth) would be the combination of frame 3 of coder 1 and frame 5 of coder 2. The procedure of determining combined average values continues until the valence value crosses again to a *non-negative* valence value. The endpoint of the audiovisual segment is then set to the frame including the offset after crossing back to a *non-negative* valence value. The ground truth of the audiovisual segment consists of the arousal and valence (A-V) values calculated.

Typically, an automatically produced segment or clip consists of a single interaction of the subject with the avatar (operator), starting with the final seconds of the avatar speaking, continuing with the subject responding (and thus reacting and expressing an emotional state audiovisually) and concluding where the avatar starts responding.

10.5.3 Feature Extraction

Our audio features include Mel-frequency Cepstrum Coefficients (MFCC) [49] and prosody features (the energy of the signal, the Root Mean Squared Energy and the pitch obtained by using a Praat pitch estimator [74]). Mel-frequency Cepstrum (MFC) is a representation of the spectrum of an audio sample which is mapped onto the nonlinear mel-scale of frequency to better approximate the human auditory system's response. The MFCC coefficients collectively make up the MFC for the specific audio segment. We used six cepstrum coefficients, thus obtaining six MFCC and six MFCC-Delta features for each audio frame. We have essentially used the typical set of features used for automatic affect recognition (e.g., [75]). Along with pitch, energy and RMS energy, we obtained a set of features with dimensionality $d = 15$ per audio frame. Note that we used a 0.04 second window with a 50% overlap (i.e. first frame 0–0.04, second from 0.02–0.06 and so on) in order to obtain a double frame rate for audio (50 Hz) compared to that of video (25 fps). This is an effective and straightforward way to synchronise the audio and video streams (similarly to [75]).



Fig. 10.4 Examples of the data at hand from the SAL database along with the extracted 20 points, used as features for the facial expression cues

To capture the facial motion displayed during a spontaneous expression we track 20 facial feature points (FFP), as illustrated in Fig. 10.4. These points are the corners of the eyebrows (4 points), eyes (8 points), nose (3 points), the mouth (4 points) and the chin (1 point). To track these facial points we used the Patras–Pantic particle filtering tracking scheme [73]. For each video segment containing n frames, we obtain a set of n vectors containing 2D coordinates of the 20 points tracked in n frames ($Tr_f = \{Tr_{f1} \dots Tr_{f20}\}$) with dimensions $n * 20 * 2$.

10.5.4 Dimensional Affect Prediction

This section describes how dimensional affect prediction from voice and face is achieved using the Bidirectional Long Short-Term Memory Neural Networks (BLSTM-NN). It first focuses on single-cue prediction from voice or face, and then introduces the model-level and output-associative fusion using the BLSTM-NNs.

10.5.4.1 Bidirectional Long Short-Term Memory Neural Networks

The traditional Recurrent Neural Networks (RNN) are unable to learn temporal dependencies longer than a few time steps due to the vanishing gradient problem [42, 43]. LSTM Neural Networks (LSTM-NNs) were introduced by Graves and Schmidhuber [32] in order to overcome this issue. The LSTM structure introduces recurrently connected memory blocks instead of traditional neural network nodes

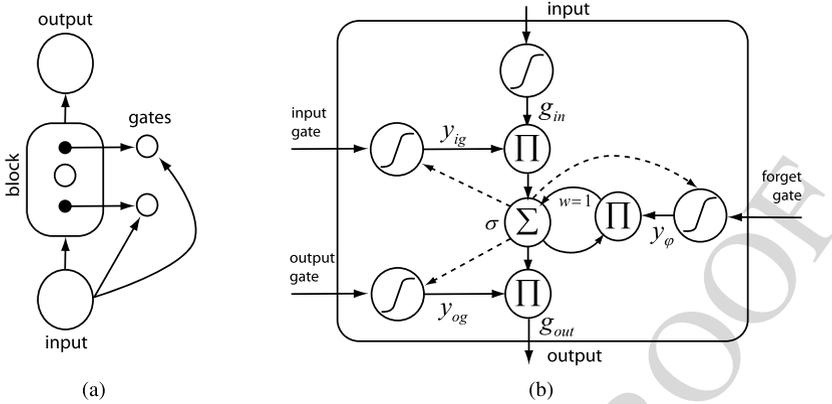


Fig. 10.5 Illustration of (a) the simplest LSTM network, with a single input, a single output, and a single memory block in place of the hidden unit, and (b) a typical implementation of an LSTM block, with multiplication units (Π), an addition unit (Σ) maintaining the cell state and typically non-linear squashing function units

(Fig. 10.5(a)). Each memory block contains memory cells and a set of multiplicative gates. In its simplest form, a memory block contains one memory cell.

As can be seen from Fig. 10.5(b), there are three types of gates: the input, output and forget gates. These gates are estimated during the training phase of an LSTM-NN.

The input, output and forget gates can be thought of as providing write, read and reset access to what is called a cell state (σ), which represents temporal network information. This can be seen from examining the state updates at time t :

$$\sigma(t) = y_{\phi}(t)\sigma(t-1) + y_{ig}(t)g_{in}(t).$$

The next state $\sigma(t)$ is defined as the sum of the forget gate at time t ($y_{\phi}(t)$) multiplied by the previous state, $\sigma(t-1)$ and the squashed input to the cell $g_{in}(t)$ multiplied by the input gate $y_{ig}(t)$. Thus, the forget gate can reset the state of the network, i.e. when $y_{\phi} \approx 0$ then the next state does not depend on the previous one:

$$\sigma(t) \approx y_{ig}(t)g_{in}(t).$$

This is similar when the input gate is near zero. Then, the next state depends only on the previous state and the forget gate value. The output of the cell is the cell state, as regulated by the value of the output gate (Fig. 10.5(b)). This configuration enforces constant error flow and overcomes the vanishing gradient problem.

In addition, traditional RNNs process input in a temporal order, thus learning input patterns by relating only to past context. Bidirectional RNNs (BRNNs) [3, 94] instead modify the learning procedure to overcome the latter issue of the past and future context: they present each of the training sequences in a forward and a backward order (to two different recurrent networks, respectively, which are connected to a common output layer). In this way, the BRNN is aware of both future and

967 past events in relation to the current timestep. The concept is directly expanded
 968 for LSTMs, referred to as Bidirectional Long Short-Term Memory neural networks
 969 (BLSTM-NN). BLSTM-NN have been shown to outperform unidirectional LSTM-
 970 NN for speech processing (e.g., [32]) and have been used for many learning tasks.
 971 They have been successfully applied to continuous emotion prediction from speech
 972 (e.g., [105, 106]) proving that modeling the sequential inputs and long range tem-
 973 poral dependencies appear to be beneficial for the task of automatic emotion predic-
 974 tion.

975 976 10.5.4.2 Single-Cue Prediction

977 The first step in continuous affect prediction task consists of prediction based on
 978 single cues. Let $\mathcal{D} = \{V, A\}$ represent the set of emotion dimensions, \mathcal{C} the set of
 979 cues consisting of the facial expressions, shoulder movement and audio cues. Given
 980 a set of input features $\mathbf{x}_c = [\mathbf{x}_{1c}, \dots, \mathbf{x}_{nc}]$ where n is the training sequence length
 981 and $c \in \mathcal{C}$, we train a machine learning technique f_d , in order to predict the relevant
 982 dimension output, $\mathbf{y}_d = [y_1, \dots, y_n]$, $d \in \mathcal{D}$.

$$983 \quad f_d : \mathbf{x} \mapsto y_d. \quad (10.2)$$

984 This step provides us with a set of predictions for each machine learning technique,
 985 and each relevant dimension employed.

986 987 10.5.4.3 Model-Level Fusion

988 As already explained in Sect. 10.4.2, since humans display multi-cue and multi-
 989 modal expressions in a complementary and redundant manner, the assumption of
 990 conditional independence between modalities and cues in decision-level fusion can
 991 result in loss of information (i.e. mutual correlation between the modalities). There-
 992 fore, we opt for model-level fusion of the continuous predictions as this has the
 993 potential of capturing correlations and structures embedded in the continuous out-
 994 put of the predictors/regressors (from different sets of cues). This is illustrated in
 995 Fig. 10.6(a).

996 More specifically, during model-level fusion, a function learns to map predictions
 997 to a dimension d from the set of cues as follows:

$$998 \quad f_{mf} : f_d(\mathbf{x}_1) \times \dots \times f_d(\mathbf{x}_m) \mapsto y_d, \quad (10.3)$$

999 where m is the total number of fused cues.

1000 1001 10.5.4.4 Output-Associative Fusion

1002 In the previous section, we have treated the prediction of valence or arousal as a
 1003 1D regression problem. However, psychological evidence shows that valence and
 1004

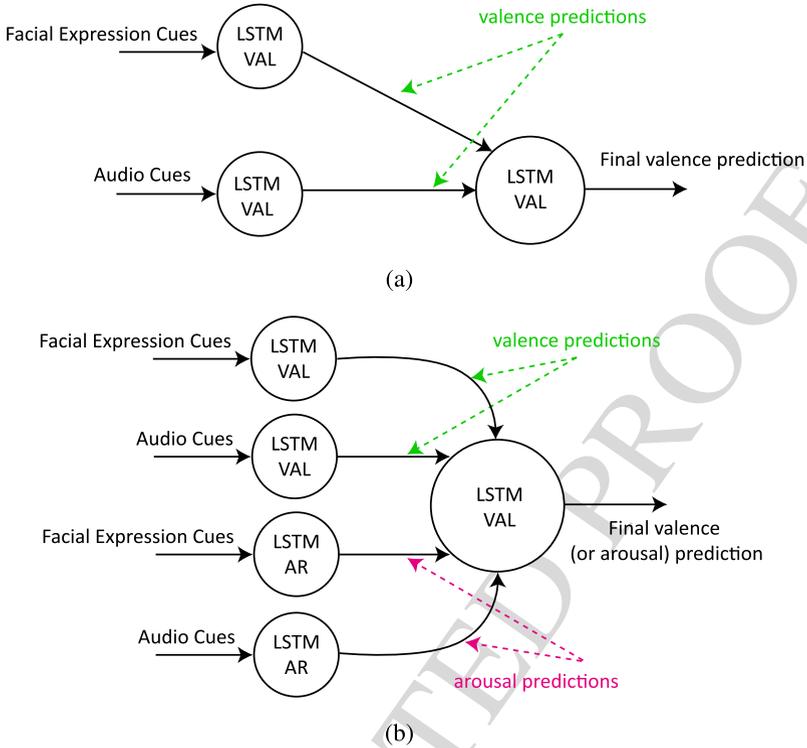


Fig. 10.6 Illustration of (a) model-level fusion and (b) output-associative fusion using facial expression and audio cues. Model-level fusion combines valence predictions from facial expression and audio cues by using a third network for the final valence prediction. Output-associative fusion combines both valence and arousal values predicted from facial expression and audio cues, again by using a third network, which outputs the final prediction.

arousal dimensions are correlated [2, 70, 107]. In order to exploit these correlations and patterns, we propose a framework capable of learning the dependencies that exist amongst the predicted dimensional values.

Given the setting described in Sect. 10.5.4.2, this framework learns to map the outputs of the intermediate predictors (each BLSTM-NN as defined in (10.2)) onto a higher (and final) level of prediction by incorporating cross-dimensional (output) dependencies (see Fig. 10.6(b)). This method, which we call *output-associative fusion*, can be represented by a function f_{oaf} :

$$f_{oaf} : f_{Ar}(\mathbf{x}_1) \times f_{Val}(\mathbf{x}_1) \times \dots \times f_{Ar}(\mathbf{x}_m) \times f_{Val}(\mathbf{x}_m) \mapsto y_d. \quad (10.4)$$

As a result, the final output, taking advantage of the temporal and bidirectional characteristics of the regressors (BLSTM-NNs), depends not only on the entire sequence of input features \mathbf{x}_i but also on the entire sequence of intermediate output predictions \mathbf{f}_d of both dimensions (see Fig. 10.6(b)).

Table 10.1 Single-cue prediction results for valence and arousal dimensions

| Dimension | Modality | RMSE | COR | SAGR |
|-----------|----------|-------|-------|-------|
| Arousal | Voice | 0.240 | 0.586 | 0.764 |
| | Face | 0.250 | 0.493 | 0.681 |
| Valence | Voice | 0.220 | 0.444 | 0.648 |
| | Face | 0.170 | 0.712 | 0.841 |

10.5.5 Experiments and Analysis

10.5.5.1 Experimental Setup

Prior to experimentation, all features have been normalized to the range of $[-1, +1]$, except for the audio features which have been found to perform better with z-normalization (i.e., normalizing to mean = 0 and standard deviation = 1).

As the main evaluation metrics we choose to use the root mean squared error (RMSE) that evaluates the root of the prediction by taking into account the squared error of the prediction from the ground truth, the correlation (COR) that provides an evaluation of the linear relationship between the prediction and the ground truth, and subsequently, an evaluation of whether the model has managed to capture linear structural patterns inhibited in the data at hand, and the sign agreement metric (SAGR) that measures the agreement level of the prediction with the ground truth by assessing the valence dimension as being positive (+) or negative (-), and the arousal dimension as being active (+) or passive (-).

For validation purposes we use a subset of the SAL-DB that consists of 134 audiovisual segments (a total of 30,042 video frames) obtained by the automatic segmentation procedure (proposed in [67]). As V-A annotations have only been provided for two female and two male subjects, for our experiments we employ *subject-dependent leave-one-sequence-out cross-validation*. More specifically, the evaluation consists of 134 folds where at each fold one sequence is left out for testing and the other 133 sequences are used for training. The prediction results are then averaged over 134 folds.

The parameter optimization for BLSTM-NNs refers to mainly determining the topology of the network along with the number of epochs, momentum and learning rate.

10.5.5.2 Results and Analysis

Single-cue results are presented in Table 10.1, while results obtained from fusion are presented in Table 10.2.

We initiate our analysis with the single-cue results (Table 10.1) and the valence dimension. Various automatic dimensional emotion prediction and recognition studies have shown that arousal can be much better predicted than valence using audio

Table 10.2 Results for output-associative fusion (AOF) and model-level fusion (MLF). The best results are obtained by employing output-associative fusion (shown in bold)

| Dimension | OAF | | | MLF | | |
|-----------|--------------|--------------|--------------|-------|-------|-------|
| | RMSE | COR | SAGR | RMSE | COR | SAGR |
| Arousal | 0.220 | 0.628 | 0.800 | 0.230 | 0.605 | 0.800 |
| Coders | 0.145 | 0.870 | 0.840 | 0.145 | 0.870 | 0.840 |
| Valence | 0.160 | 0.760 | 0.892 | 0.170 | 0.748 | 0.856 |
| Coders | 0.141 | 0.850 | 0.860 | 0.141 | 0.850 | 0.860 |

cues (e.g., [33, 68, 100, 105]). Our experimental results also support these findings indicating that the visual cues appear more informative for predicting the valence dimension. The facial expression cues provide a higher correlation with the ground truth (COR = 0.71) compared to the audio cues (COR = 0.44). This fact is also confirmed by the RMSE and SAGR metrics. The facial expression cues also provide higher SAGR (0.84), indicating that the predictor was accurate in predicting an emotional state as positive or negative for 84% of the frames. For prediction of the arousal dimension the audio cues appear to be superior to the visual cues. More specifically, audio cues provide COR = 0.59, whereas the facial expression cues provide COR = 0.49.

Fusing facial and audio cues using model-level fusion outperforms the single-cue prediction results. Model-level fusion appears to be much better for predicting the valence dimension rather than the arousal dimension. This is mainly due to the fact that the single-cue predictors for valence dimension perform better, thus containing more correct temporal dependencies and structural characteristics (while the weaker arousal predictors contain fewer of these dependencies). Model-level fusion also re-confirms that visual cues are more informative for valence dimension than the audio cues. Finally, the newly proposed output-associative fusion provides the best results, outperforming both single-cue analysis and model-level fusion results. We denote that the performance increase of output-associative fusion is higher for the arousal dimension (compared to the valence dimension). This could be justified by the fact that the single-cue predictors for valence perform better than for arousal (Table 10.1) and thus, more correct valence patterns are passed onto the output-associative fusion framework. An example of the output-associative valence and arousal prediction from face and audio is shown in Fig. 10.7.

Based on the experimental results provided in Tables 10.1–10.2, we conclude the following.

- Facial expression cues are better suited to the task of continuous valence prediction compared to audio cues. For arousal dimension, instead, the audio cues appear to perform better. This is in accordance with the previous findings in the literature.
- The inherent temporal and structured nature of continuous affective data appears to be highly suitable for predictors that can model temporal dependencies and relate temporally distant events. To evaluate the performance of such frameworks,

1151 the use of not only the RMSE but also the correlation coefficient appears to be
1152 very important. Furthermore, the use of other emotion-specific metrics, such as
1153 the SAGR (used in this work), is also desirable as they contain valuable information
1154 regarding emotion-specific aspects of the predictions.

- 1155 • As confirmed by the psychological theory, valence and arousal are correlated.
1156 Such correlations appear to exist in our data where fusing predictions from both
1157 valence and arousal dimensions (output-associative fusion) improves the results
1158 compared to using predictions from either valence or arousal dimension alone (as
1159 in the model-level fusion case).
- 1160 • In general, audiovisual data appear to be more useful for predicting valence than
1161 for predicting arousal. While arousal is better predicted by using audio features
1162 alone, valence is better predicted by using audiovisual data.

1163 Overall, our output-associative fusion framework (i) achieves $RMSE = 0.160$,
1164 $COR \approx 0.760$ and $SAGR \approx 0.900$ for the valence dimension, compared to the human
1165 coder (inter-coder) $RMSE \approx 0.141$, $COR \approx 0.850$, and $SAGR \approx 0.860$, and
1166 (ii) provides $RMSE = 0.220$, $COR \approx 0.628$ and $SAGR \approx 0.800$ for the arousal
1167 dimension, compared to the human coder (inter-coder) $RMSE \approx 0.145$, $COR \approx 0.870$
1168 and $SAGR \approx 0.840$.

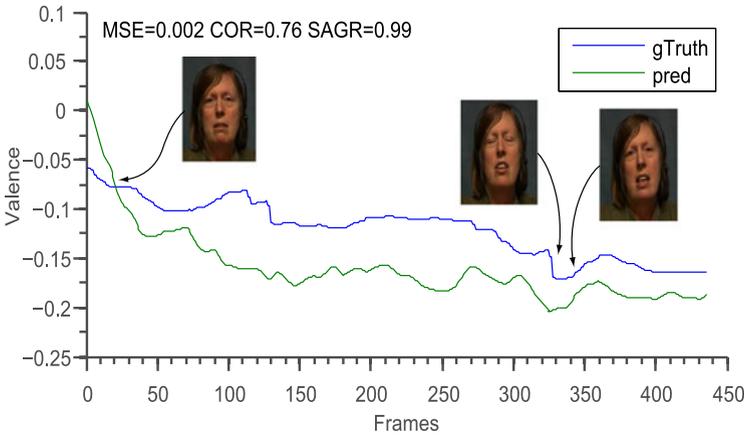
1169 In our experiments we employed a subject-dependent leave-one-sequence-out
1170 cross-validation procedure due to the small number of annotated data available. As
1171 spontaneous expressions appear to have somewhat person-dependent characteristics,
1172 subject-independent experimentation is likely to be more challenging and affect
1173 our prediction results.

1174 1175 1176 1177 **10.6 Concluding Remarks**

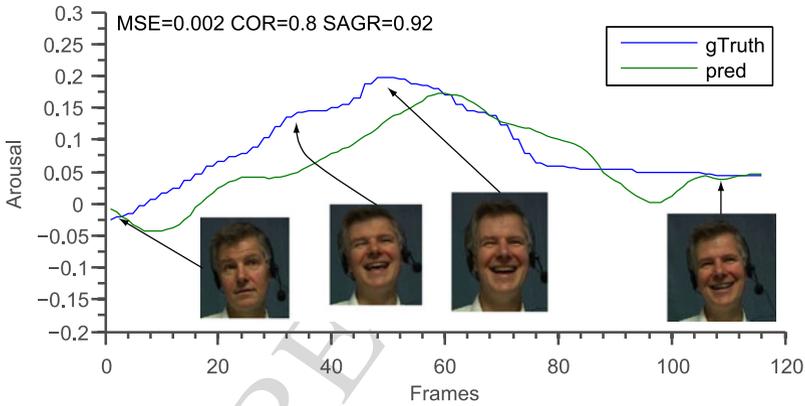
1178 The review provided in this chapter suggests that the automatic affect sensing field
1179 has slowly started shifting from categorical (and discrete) affect recognition to di-
1180 mensional (and continuous) affect prediction to be able to capture the complexity of
1181 affect expressed in naturalistic settings. There is a growing research interest driven
1182 by various advances and demands (e.g., real-time representation and analysis of
1183 naturalistic and continuous human affective behavior for emotion-related disorders
1184 like autism), and funded by various research projects (e.g., European Union FP 7,
1185 SEMAINE¹). To date, despite the existence of a number of dimensional emotion
1186 models, the two-dimensional model of arousal and valence appears to be the most
1187 widely used model in automatic measurement of affect from audio, visual and bio
1188 signals.

1189 The current automatic measurement technology has already started dealing with
1190 spontaneous data obtained in less-controlled environments using various sensing
1191 devices, and exploring a number of machine learning techniques and evaluation
1192 measures. However, naturalistic settings pose many challenges to continuous affect
1193

1194
1195 ¹<http://www.semaine-project.eu>
1196



(a)



(b)

Fig. 10.7 Valence and arousal ground truth (gTruth) compared to predictions (pred) from output-associative fusion of facial expressions and audio cues

sensing and prediction (e.g., when subjects are not restricted in terms of mobility, the level of noise in all recorded signals tends to increase), as well as affect synthesis and generation. As a consequence, a number of issues that should be addressed in order to advance the field remain unclear. These have been summarized and discussed in this chapter.

As summarized in Sect. 10.4.2 and reviewed in [37], to date, only a few systems have actually achieved dimensional affect prediction from multiple modalities. Overall, existing systems use different training/testing datasets (which differ in the way affect is elicited and annotated), they differ in the underlying affect model (i.e., target affect categories), as well as in the employed modality or combination of

modalities, and the applied evaluation method. As a consequence, it remains unclear which recognition and prediction method is suitable for dimensional affect prediction from which modalities and cues. These challenges should be addressed in order to advance the field while identifying the importance, as well as the feasibility, of the following issues:

1. *Among the available remotely observable and remotely unobservable modalities, which ones should be used for automatic dimensional affect prediction? Should we investigate the innate priority among the modalities to be preferred for each affect dimension? Does this depend on the context (who the subject is, where she is, what her current task is, and when the observed behavior has been shown)?*

Continuous long-term monitoring of bio signals (e.g., autonomic nervous system) appears to be particularly useful and usable for health care applications (e.g., stress and pain monitoring, autism-related assistive technology). Using bio signals for automatic measurement is especially important for applications where people do not easily express themselves outwardly with facial and bodily expressions (e.g., people with autism spectrum disorders) [24]. As stated before, various automatic dimensional emotion prediction and recognition studies have shown that arousal can be much better predicted than valence using audio cues (e.g., [33, 68, 100, 105]). For the valence dimension instead, visual cues (e.g., facial expressions and shoulder movements) appear to perform better [68]. Whether such conclusions hold for different contexts and different data remains to be evaluated. Another significant research finding is that when multiple modalities are available during data annotation, both speed and accuracy of judgments increase when the modalities are expressing the same emotion [15]. How such findings should be incorporated into automatic dimensional affect predictors remains to be researched further.

2. *When labeling emotions, which signals better correlate with self assessment and which ones correlate with independent observer assessment?*

When acquiring and annotating emotional data, there exist individual differences in emotional response, as well as individual differences in the use of rating scales. We have mentioned some of these differences before, in Sect. 10.4.1. Research also shows that affective state labeling is significantly affected by factors such as familiarity of the person and context of the interaction [44]. Even if the emotive patterns to be labeled are fairly similar, human perception is biased by context and prior experience. Moreover, Feldman presented evidence that when individuals are shown emotional stimulus, they differ in their attention to valence and arousal dimensions [23]. We have also mentioned cross-cultural intensity differences in labeling emotional behaviors [96]. If such issues are ignored and the ratings provided by the human annotators are simply averaged, the measure obtained may be useful in certain experimental contexts but it will be insensitive to individual variations in subjective experience. More specifically, this will imply having a scale that assumes that individual differences are unimportant or nonexistent. An implication of this view is that for an ideal representation of a subject's affective state, labeling schemes and rating scales should be clearly defined (e.g., by making the subjective distances between adjacent numbers on every portion

of the scale equal) and contextualized (e.g., holding the environmental cues constant), both self assessment and external observer assessment (preferably from observers who are familiar with the user to be assessed) should be obtained and used, and culture-related issues should be taken into consideration.

3. *How does the baseline problem affect prediction? Is an objective basis (e.g., a frame with an expressionless display) strictly needed prior to computing the dimensional affect values? If so, how can this be obtained in a fully automatic manner from naturalistic data?*

Determining the baseline in naturalistic affective displays is challenging even for human observers. This is particularly the case for the visual modality which constitutes of varying head pose and head gestures (like nods and shakes), speech-related facial actions, and blended facial expressions. The implications for automatic analysis can initially be addressed by training predictors that predict baseline (or neutrality) for each cue and modality separately.

4. *How should intensity be modeled for dimensional and continuous affect prediction? Should the aim be personalizing systems for each subject, or creating systems that are expected to generalize across subjects?*

Modeling the intensity of emotions should be based on the task-dependent environment and target user group. A common way of measuring affect from bio signals is doing it for each participant separately (without computing baseline), e.g., [10]. Similarly to the recent works on automatic affect prediction from the audio or the visual cues (e.g., [69]), better insight may be obtained by comparing subject-dependent vs. subject-independent prediction results. Customizing the automatic predictors to specific user needs is usually desired and advantageous.

5. *In a continuous affect space, how should duration of affect be defined? How can this be incorporated in automated systems? Will focusing on shorter or longer observations affect the accuracy of the measurement process?*

Similarly to modeling the emotional intensity level, determining the affect duration should be done based on the task-dependent environment and target user group. Focusing on shorter or longer durations appears to have an effect on the prediction accuracy. Achieving real-time affect prediction requires a small window size to be used for analysis (i.e., a few seconds, e.g., [10]), while on the other hand obtaining a reliable prediction accuracy requires long(er)-term monitoring [6, 83]. Therefore, analysis duration should be determined as a trade-off between reliable prediction accuracy and real-time requirements of the automatic system.

Finding comprehensive and thorough answers to the questions posed above, and fully exploring the terrain of the dimensional and continuous affect prediction, depends on all relevant research fields (engineering, computer science, psychology, neuroscience, and cognitive sciences) stepping out of their labs, working side-by-side together on real-life applications, and sharing the experience and the insight acquired on the way, to make affect research tangible for realistic settings and lay people [76]. Pioneering projects representing such inter-disciplinary efforts have already started emerging, ranging, for instance, from publishing compiled books of related work (e.g., [30]) and organizing emotion recognition challenges (e.g., INTERSPEECH 2010 Paralinguistic Challenge featuring the affect sub-challenge

with a focus on dimensional affect [93]) to projects as varied as affective human-embodied conversational agent interaction (e.g., European Union FP 7 SEMAINE [89, 90]), and affect sensing for autism (e.g., [76, 78]).

10.7 Summary

Human affective behavior is multimodal, continuous and complex. Despite major advances within the affective computing research field, modeling, analyzing, interpreting and responding to human affective behavior still remains a challenge for automated systems as affect and emotions are complex constructs, with fuzzy boundaries and with substantial individual differences in expression and experience [7]. Therefore, affective and behavioral computing researchers have recently invested increased effort in exploring how to best model, analyze and interpret the subtlety, complexity and continuity (represented along a continuum e.g., from -1 to $+1$) of affective behavior in terms of latent dimensions (e.g., arousal, power and valence) and appraisals, rather than in terms of a small number of discrete emotion categories (e.g., happiness and sadness). This chapter aimed to (i) give a brief overview of the existing efforts and the major accomplishments in modeling and analysis of emotional expressions in dimensional and continuous space while focusing on open issues and new challenges in the field, and (ii) introduce a representative approach for multimodal continuous analysis of affect from voice and face, and provide experimental results using the audiovisual Sensitive Artificial Listener (SAL) Database of natural interactions. The chapter concluded by posing a number of questions that highlight the significant issues in the field, and by extracting potential answers to these questions from the relevant literature.

10.8 Questions

1. What are the major approaches used for affect modeling and representation? How do they differ from each other?
2. Why has the dimensional affect representation gained interest?
3. What are the dimensions used for representing emotions?
4. Affect research scientists usually make a number of assumptions and simplifications while studying emotions. What are these assumptions and simplifications? What implications do they have?
5. How is human affect sensed and measured? What are the signals measured for analyzing human affect?
6. How are affective data acquired and annotated?
7. What is the current state of the art in automatic affect prediction and recognition?
8. What are the challenges faced in automatic dimensional affect recognition?
9. List a number of applications that use the dimensional representation of emotions.

- 1381 10. What features are extracted to represent an audio-visual affective sequence?
1382 How are the audio and video streams synchronized?
- 1383 11. What is a Bidirectional Long Short-Term Memory Neural Network? How does
1384 it differ from a traditional Recurrent Neural Network?
- 1385 12. What is meant by the statement ‘valence and arousal dimensions are corre-
1386 lated’? What implications does this have on automatic affect prediction?
- 1387 13. What is output-associative fusion? How does it compare to model-level fusion?
- 1388 14. How are the root mean squared error, correlation, and sign agreement used for
1389 evaluating the automatic prediction of emotions?
1390

1391 10.9 Glossary

- 1394 • *Categorical description of affect.* Hypothesizes that there exist a small number
1395 of emotions categories (i.e., anger, disgust, fear, happiness, sadness and surprise)
1396 that are basic, hard-wired in our brain, and recognized universally (e.g. [18]).
- 1397 • *Dimensional description of affect.* Hypothesizes that affective states are not inde-
1398 pendent from one another; rather, they are related to one another in a systematic
1399 manner.
- 1400 • *Circumplex Model of Affect.* A circular configuration introduced by Russell [82],
1401 based on the hypothesis that each basic emotion represents a bipolar entity being
1402 a part of the same emotional continuum.
- 1403 • *PAD emotion space.* The three dimensional description of emotion in terms of
1404 pleasure–displeasure, arousal–nonarousal and dominance–submissiveness [63].
- 1405 • *Dimensional and continuous affect prediction.* Analyzing and inferring the sub-
1406 tlety, complexity and continuity of affective behavior in terms of latent dimen-
1407 sions (e.g., valence and arousal) by representing it along a continuum (e.g., from
1408 -1 to $+1$) without discretization.
- 1409 • *Long Short-Term Memory neural network.* A Bidirectional Recurrent Neural Net-
1410 work that consists of recurrently connected memory blocks, and uses input, out-
1411 put and forget gates to represent and learn the temporal information and depen-
1412 dencies.
- 1413 • *Output-associative fusion.* A fusion approach that uses multi-layered prediction,
1414 i.e. the initial features extracted from each modality are used for intermediate
1415 (output) prediction, and these are further used for a higher (and final) level of
1416 prediction (by incorporating cross-dimensional dependencies).
1417

1418 **Acknowledgements** This work has been funded by EU [FP7/2007-2013] Grant agreement
1419 No. 211486 (SEMAINE) and the ERC Starting Grant agreement No. ERC-2007-StG-203143
1420 (MAHNOB).
1421

1422 10.10 References

- 1423 1. Affectiva’s homepage: <http://www.affectiva.com/> (2011)
1424
1425
1426

- 1427 2. Alvarado, N.: Arousal and valence in the direct scaling of emotional response to film clips.
1428 *Motiv. Emot.* **21**, 323–348 (1997)
- 1429 3. Baldi, P., Brunak, S., Frasconi, P., Pollastri, G., Soda, G.: Exploiting the past and the future
1430 in protein secondary structure prediction. *Bioinformatics* **15**, 937–946 (1999)
- 1431 4. Bartneck, C.: Integrating the occ model of emotions in embodied characters. In: Proc. of the
1432 Workshop on Virtual Conversational Characters, pp. 39–48 (2002)
- 1433 5. Beck, A., Canamero, L., Bard, K.A.: Towards an affect space for robots to display emotional
1434 body language. In: Proc. IEEE Int. Symp. in Robot and Human Interactive Communication,
1435 pp. 464–469 (2010)
- 1436 6. Berntson, G.G., Bigger, J.T., Eckberg, D.L., Grossman, P., Kaufmann, P.G., Malik, M., Na-
1437 garaja, H.N., Porges, S.W., Saul, J.P., Stone, P.H., van der Molen, M.W.: Heart rate variabil-
1438 ity: origins, methods, and interpretive caveats. *Psychophysiology* **34**(6), 623 (1997)
- 1439 7. Calvo, R.A., D’Mello, S.: Affect detection: An interdisciplinary review of models, methods,
1440 and their applications. *IEEE Trans. Affect. Comput.* **1**(1), 18–37 (2010)
- 1441 8. Caridakis, G., Malatesta, L., Kessous, L., Amir, N., Raouzaoui, A., Karpouzis, K.: Modeling
1442 naturalistic affective states via facial and vocal expressions recognition. In: Proc. of ACM
1443 Int. Conf. on Multimodal Interfaces, pp. 146–154 (2006)
- 1444 9. Chanel, G., Ansari-Asl, K., Pun, T.: Valence-arousal evaluation using physiological signals in
1445 an emotion recall paradigm. In: Proc. of IEEE Int. Conf. on Systems, Man and Cybernetics,
1446 pp. 2662–2667, October 2007
- 1447 10. Chanel, G., Kierkels, J.J.M., Soleymani, M., Pun, T.: Short-term emotion assessment in a
1448 recall paradigm. *Int. J. Hum.-Comput. Stud.* **67**(8), 607–627 (2009)
- 1449 11. Cowie, R., Douglas-Cowie, E., Savvidou, S., McMahon, E., Sawey, M., Schroder, M.: Feel-
1450 trace: An instrument for recording perceived emotion in real time. In: Proc. of ISCA Work-
1451 shop on Speech and Emotion, pp. 19–24 (2000)
- 1452 12. Cowie, R., Douglas-Cowie, E., Tsapatsoulis, N., Votsis, G., Kollias, S., Fellenz, W., Taylor,
1453 J.G.: Emotion recognition in human–computer interaction. *IEEE Signal Process. Mag.* **18**,
1454 33–80 (2001)
- 1455 13. Cowie, R., Gunes, H., McKeown, G., Vaclau-Schneider, L., Armstrong, J., Douglas-Cowie,
1456 E.: The emotional and communicative significance of head nods and shakes in a naturalistic
1457 database. In: Proc. of LREC Int. Workshop on Emotion, pp. 42–46 (2010)
- 1458 14. Davitz, J.: Auditory correlates of vocal expression of emotional feeling. In: *The Communi-
1459 cation of Emotional Meaning*, pp. 101–112. McGraw-Hill, New York (1964)
- 1460 15. de Gelder, B., Vroomen, J.: The perception of emotions by ear and by eye. *Cogn. Emot.* **23**,
1461 289–311 (2000)
- 1462 16. Douglas-Cowie, E., Cowie, R., Sneddon, I., Cox, C., Lowry, L., McRorie, M., Martin, L.
1463 Jean-Claude, Devillers, J.-C., Abrilian, A., Batliner, S., Noam, A., Karpouzis, K.: The HU-
1464 MAINE database: addressing the needs of the affective computing community. In: Proc.
1465 of the Second Int. Conf. on Affective Computing and Intelligent Interaction, pp. 488–500
1466 (2007)
- 1467 17. Ekman, P., Friesen, W.V.: Head and body cues in the judgment of emotion: A reformulation.
1468 *Percept. Mot. Skills* **24**, 711–724 (1967)
- 1469 18. Ekman, P., Friesen, W.V.: *Unmasking the Face: A Guide to Recognizing Emotions from
1470 Facial Clues*. Prentice Hall, New Jersey (1975)
- 1471 19. Espinosa, H.P., Garcia, C.A.R., Pineda, L.V.: Features selection for primitives estimation on
1472 emotional speech. In: Proc. IEEE Int. Conf. on Acoustics Speech and Signal Processing,
pp. 5138–5141 (2010)
20. Eyben, F., Wöllmer, M., Poitschke, T., Schuller, B., Blaschke, C., Färber, B., Nguyen-
Thien, N.: Emotion on the road—necessity, acceptance, and feasibility of affective com-
puting in the car. *Adv. Hum.-Comput. Interact.* **2010**, 263593 (2010), 17 pages
21. Eyben, F., Wöllmer, M., Valstar, M., Gunes, H., Schuller, B., Pantic, M.: String-based au-
diovisual fusion of behavioural events for the assessment of dimensional affect. In: Proc. of
IEEE Int. Conf. on Automatic Face and Gesture Recognition (2011)

- 1473 22. Faghihi, U., Fournier-Viger, P., Nkambou, R., Poirier, P., Mayers, A.: How emotional mechanism helps episodic learning in a cognitive agent. In: Proc. IEEE Symp. on Intelligent Agents, pp. 23–30 (2009)
- 1474
- 1475 23. Feldman, L.: Valence focus and arousal focus: Individual differences in the structure of affective experience. *J. Pers. Soc. Psychol.* **69**, 153–166 (1995)
- 1476
- 1477 24. Fletcher, R., Dobson, K., Goodwin, M.S., Eydgahi, H., Wilder-Smith, O., Fernholz, D., Kuboyama, Y., Hedman, E., Poh, M.Z., Picard, R.W.: iCalm: Wearable sensor and network architecture for wirelessly communicating and logging autonomic activity. *IEEE Tran. on Information Technology in Biomedicine* **14**(2), 215
- 1478
- 1479
- 1480 25. Fontaine, J.R., Scherer, K.R., Roesch, E.B., Ellsworth, P.: The world of emotion is not two-dimensional. *Psychol. Sci.* **18**, 1050–1057 (2007)
- 1481
- 1482 26. Fragopanagos, N., Taylor, J.G.: Emotion recognition in human–computer interaction. *Neural Netw.* **18**(4), 389–405 (2005)
- 1483
- 1484 27. Frijda, N.H.: *The Emotions*. Cambridge University Press, Cambridge (1986)
- 1485
- 1486 28. Gilroy, S.W., Cavazza, M., Niiranen, M., Andre, E., Vogt, T., Urbain, J., Benayoun, M., Seichter, H., Billingham, M.: Pad-based multimodal affective fusion. In: Proc. Int. Conf. on Affective Computing and Intelligent Interaction Workshops, pp. 1–8 (2009)
- 1487
- 1488 29. Glowinski, D., Camurri, A., Volpe, G., Dael, N., Scherer, K.: Technique for automatic emotion recognition by body gesture analysis. In: Proc. of Computer Vision and Pattern Recognition Workshops, pp. 1–6 (2008)
- 1489
- 1490 30. Gokcay, D., Yildirim, G.: *Affective Computing and Interaction: Psychological, Cognitive and Neuroscientific Perspectives*. IGI Global, Hershey (2011)
- 1491
- 1492 31. Grandjean, D., Sander, D., Scherer, K.R.: Conscious emotional experience emerges as a function of multilevel, appraisal-driven response synchronization. *Conscious. Cogn.* **17**(2), 484–495 (2008)
- 1493
- 1494 32. Graves, A., Schmidhuber, J.: Framewise phoneme classification with bidirectional LSTM and other neural network architectures. *Neural Netw.* **18**, 602–610 (2005)
- 1495
- 1496 33. Grimm, M., Kroschel, K.: Emotion estimation in speech using a 3d emotion space concept. In: Proc. IEEE Automatic Speech Recognition and Understanding Workshop, pp. 381–385 (2005)
- 1497
- 1498 34. Grimm, M., Mower, E., Kroschel, K., Narayanan, S.: Primitives based estimation and evaluation of emotions in speech. *Speech Commun.* **49**, 787–800 (2007)
- 1499
- 1500 35. Grimm, M., Kroschel, K., Narayanan, S.: The Vera am Mittag German audio-visual emotional speech database. In: ICME, pp. 865–868. IEEE Press, New York (2008)
- 1501
- 1502 36. Grundlehner, B., Brown, L., Penders, J., Gyselinckx, B.: The design and analysis of a real-time, continuous arousal monitor. In: Proc. Int. Workshop on Wearable and Implantable Body Sensor Networks, pp. 156–161 (2009)
- 1503
- 1504 37. Gunes, H., Pantic, M.: Automatic, dimensional and continuous emotion recognition. *Int. J. Synth. Emot.* **1**(1), 68–99 (2010)
- 1505
- 1506 38. Gunes, H., Pantic, M.: Automatic measurement of affect in dimensional and continuous spaces: Why, what, and how. In: Proc. of Measuring Behavior, pp. 122–126 (2010)
- 1507
- 1508 39. Gunes, H., Pantic, M.: Dimensional emotion prediction from spontaneous head gestures for interaction with sensitive artificial listeners. In: Proc. of International Conference on Intelligent Virtual Agents, pp. 371–377 (2010)
- 1509
- 1510 40. Gunes, H., Piccardi, M., Pantic, M.: Affective computing: focus on emotion expression, synthesis, and recognition. In: Or, J. (ed.) *From the Lab to the Real World: Affect Recognition using Multiple Cues and Modalities*, pp. 185–218. I-Tech Education and Publishing, Vienna (2008)
- 1511
- 1512
- 1513 41. Haag, A., Goronzy, S., Schaich, P., Williams, J.: Emotion recognition using bio-sensors: First steps towards an automatic system. In: LNCS, vol. 3068, pp. 36–48 (2004)
- 1514
- 1515 42. Hochreiter, S.: *Untersuchungen zu dynamischen neuronalen Netzen*. Diploma thesis, Institut für Informatik, Lehrstuhl Prof. Brauer, Technische Universität München (1991)
- 1516
- 1517 43. Hochreiter, S.: The vanishing gradient problem during learning recurrent neural nets and problem solutions. *Int. J. Uncertain. Fuzziness Knowl.-Based Syst.* **6**(2), 107–116 (1998)
- 1518

- 1519 44. Hoque, M.E., El Kaliouby, R., Picard, R.W.: When human coders (and machines) disagree
 1520 on the meaning of facial affect in spontaneous videos. In: Proc. of Intelligent Virtual Agents,
 1521 pp. 337–343 (2009)
- 1522 45. Huang, T.S., Hasegawa-Johnson, M.A., Chu, S.M., Zeng, Z., Tang, H.: Sensitive talking
 1523 heads. *IEEE Signal Process. Mag.* **26**, 67–72 (2009)
- 1524 46. Hutter, G.L.: Relations between prosodic variables and emotions in normal American english
 1525 utterances. *J. Speech Hear. Res.* **11**, 481–487 (1968)
- 1526 47. Ioannou, S., Raouzaïou, A., Tzouvaras, V., Mailis, T., Karpouzis, K., Kollias, S.: Emotion
 1527 recognition through facial expression analysis based on a neurofuzzy method. *Neural Netw.*
 1528 **18**(4), 423–435 (2005)
- 1529 48. Jia, J., Zhang, S., Meng, F., Wang, Y., Cai, L.: Emotional audio-visual speech synthesis based
 1530 on PAD. *IEEE Trans. Audio Speech Lang. Process.* **PP**(9), 1 (2010)
- 1531 49. Jurafsky, D., Martin, J.H.: *Speech and Language Processing: An Introduction to Natural Lan-
 1532 guage Processing, Computational Linguistics and Speech Recognition*, 2nd edn. Prentice-
 1533 Hall, New York (2008)
- 1534 50. Kanluan, I., Grimm, M., Kroschel, K.: Audio-visual emotion recognition using an emotion
 1535 recognition space concept. In: Proc. of the 16th European Signal Processing Conference
 1536 (2008)
- 1537 51. Karg, M., Schwimmbeck, M., Kühnlenz, K., Buss, M.: Towards mapping emotive gait pat-
 1538 terns from human to robot. In: Proc. IEEE Int. Symp. in Robot and Human Interactive Com-
 1539 munication, pp. 258–263 (2010)
- 1540 52. Khalili, Z., Moradi, M.H.: Emotion recognition system using brain and peripheral signals:
 1541 Using correlation dimension to improve the results of EEG. In: Proc. Int. Joint Conf. on
 1542 Neural Networks, pp. 1571–1575 (2009)
- 1543 53. Kierkels, J.J.M., Soleymani, M., Pun, T.: Queries and tags in affect-based multimedia re-
 1544 trieval. In: Proc. IEEE Int. Conf. on Multimedia and Expo, pp. 1436–1439 (2009)
- 1545 54. Kim, J.: Robust speech recognition and understanding. In: Grimm, M., Kroschel, K. (eds.)
 1546 *Bimodal Emotion Recognition using Speech and Physiological Changes*, pp. 265–280.
 1547 I-Tech Education and Publishing, Vienna (2007)
- 1548 55. Kim, J., Andre, E.: Emotion recognition based on physiological changes in music listening.
 1549 *IEEE Trans. Pattern Anal. Mach. Intell.* **30**(12), 2067–2083 (2008)
- 1550 56. Kipp, M., Martin, J.-C.: Gesture and emotion: Can basic gestural form features discriminate
 1551 emotions? In: Proc. Int. Conf. on Affective Computing and Intelligent Interaction Work-
 1552 shops, pp. 1–8 (2009)
- 1553 57. Kleinsmith, A., Bianchi-Berthouze, N.: Recognizing affective dimensions from body pos-
 1554 ture. In: Proc. of the Int. Conf. on Affective Computing and Intelligent Interaction, pp. 48–58
 1555 (2007)
- 1556 58. Kleinsmith, A., De Silva, P.R., Bianchi-Berthouze, N.: Recognizing emotion from postures:
 1557 Cross-cultural differences in user modeling. In: Proc. of the Conf. on User Modeling, pp. 50–
 1558 59 (2005)
- 1559 59. Kulic, D., Croft, E.A.: Affective state estimation for human-robot interaction. *IEEE Trans.*
 1560 *Robot.* **23**(5), 991–1000 (2007)
- 1561 60. Lang, P.J.: *The Cognitive Psychophysiology of Emotion: Anxiety and the Anxiety Disorders*.
 1562 Erlbaum, Hillside (1985)
- 1563 61. Levenson, R.: Emotion and the autonomic nervous system: A prospectus for research on
 1564 autonomic specificity. In: *Social Psychophysiology and Emotion: Theory and Clinical Ap-
 plications*, pp. 17–42 (1988)
62. McKeown, G., Valstar, M.F., Cowie, R., Pantic, M.: The SEMAINE corpus of emotionally
 coloured character interactions. In: Proc. of IEEE Int'l Conf. Multimedia, Expo (ICME'10),
 pp. 1079–1084, July 2010
63. Mehrabian, A.: Pleasure-arousal-dominance: A general framework for describing and mea-
 suring individual differences in temperament. *Curr. Psychol.* **14**, 261–292 (1996)
64. Mihelj, M., Novak, D., Muni, M.: Emotion-aware system for upper extremity rehabilitation.
 In: Proc. Int. Conf. on Virtual Rehabilitation, pp. 160–165 (2009)

- 1565 65. Nakasone, A., Prendinger, H., Ishizuka, M.: Emotion recognition from electromyography
1566 and skin conductance. In: Proc. of the 5th International Workshop on Biosignal Interpretation,
1567 pp. 219–222 (2005)
- 1568 66. Nicolaou, M.A., Gunes, H., Pantic, M.: Audio-visual classification and fusion of sponta-
1569 neous affective data in likelihood space. In: Proc. of IEEE Int. Conf. on Pattern Recognition,
1570 pp. 3695–3699 (2010)
- 1571 67. Nicolaou, M.A., Gunes, H., Pantic, M.: Automatic segmentation of spontaneous data using
1572 dimensional labels from multiple coders. In: Proc. of LREC Int. Workshop on Multimodal
1573 Corpora: Advances in Capturing, Coding and Analyzing Multimodality, pp. 43–48 (2010)
- 1574 68. Nicolaou, M.A., Gunes, H., Pantic, M.: Continuous prediction of spontaneous affect from
1575 multiple cues and modalities in valence–arousal space. *IEEE Trans. Affect. Comput.* **2**(2),
1576 92–105 (2011)
- 1577 69. Nicolaou, M.A., Gunes, H., Pantic, M.: Output-associative RVM regression for dimensional
1578 and continuous emotion prediction. In: Proc. of IEEE Int. Conf. on Automatic Face and
1579 Gesture Recognition (2011)
- 1580 70. Oliveira, A.M., Teixeira, M.P., Fonseca, I.B., Oliveira, M.: Joint model-parameter validation
1581 of self-estimates of valence and arousal: Probing a differential-weighting model of affective
1582 intensity. In: Proc. of the 22nd Annual Meeting of the Int. Society for Psychophysics,
1583 pp. 245–250 (2006)
- 1584 71. Ortony, A., Clore, G.L., Collins, A.: *The Cognitive Structure of Emotions*. Cambridge Uni-
1585 versity Press, Cambridge (1988)
- 1586 72. Parkinson, B.: *Ideas and Realities of Emotion*. Routledge, London (1995)
- 1587 73. Patras, I., Pantic, M.: Particle filtering with factorized likelihoods for tracking facial features.
1588 In: Proc. of the IEEE Int. Conf. on Automatic Face and Gesture Recognition, pp. 97–102
1589 (2004)
- 1590 74. Paul, B.: Accurate short-term analysis of the fundamental frequency and the harmonics-to-
1591 noise ratio of a sampled sound. In: Proceedings of the Institute of Phonetic Sciences, pp. 97–
1592 110 (1993)
- 1593 75. Petridis, S., Gunes, H., Kaltwang, S., Pantic, M.: Static vs. dynamic modeling of human
1594 nonverbal behavior from multiple cues and modalities. In: Proc. of ACM Int. Conf. on Mul-
1595 timodal Interfaces, pp. 23–30 (2009)
- 1596 76. Picard, R.W.: Emotion research by the people, for the people. *Emotion Review* **2**(3), 250–
1597 254
- 1598 77. Plutchik, R., Conte, H.R.: *Circumplex Models of Personality and Emotions*. APA, Washing-
1599 ton (1997)
- 1600 78. Poh, M.Z., Swenson, N.C., Picard, R.W.: A wearable sensor for unobtrusive, long-term as-
1601 sessment of electrodermal activity. *IEEE Trans. Inf. Technol. Biomed.* **57**(5), 1243–1252
1602 (2010)
- 1603 79. Pun, T., Alecu, T.I., Chanel, G., Kronegg, J., Voloshynovskiy, S.: Brain–computer interaction
1604 research at the Computer Vision and Multimedia Laboratory, University of Geneva. *IEEE*
1605 *Trans. Neural Syst. Rehabil. Eng.* **14**, 210–213 (2006)
- 1606 80. Rehm, M., Wissner, M.: Gamble-a multiuser game with an embodied conversational agent.
1607 In: *Lecture Notes in Computer Science*, vol. 3711, pp. 180–191 (2005)
- 1608 81. Roseman, I.J.: Cognitive determinants of emotion: A structural theory. In: Shaver, P. (ed.)
1609 *Review of Personality & Social Psychology*, Beverly Hills, CA, vol. 5, pp. 11–36. Sage,
1610 Thousand Oaks (1984)
- 1610 82. Russell, J.A.: A circumplex model of affect. *J. Pers. Soc. Psychol.* **39**, 1161–1178 (1980)
- 1610 83. Salahuddin, L., Cho, J., Jeong, M.G., Kim, D.: Ultra short term analysis of heart rate variabil-
1610 ity for monitoring mental stress in mobile settings. In: Proc. of the IEEE 29th International
1610 Conference of the EMBS, pp. 39–48 (2007)
- 1610 84. Sander, D., Grandjean, D., Scherer, K.R.: A systems approach to appraisal mechanisms in
1610 emotion. *Neural Netw.* **18**(4), 317–352 (2005)
- 1610 85. Scherer, K.R., Oshinsky, J.S.: Cue utilization in emotion attribution from auditory stimuli.
1610 *Motiv. Emot.* **1**, 331–346 (1977)

- 1611 86. Scherer, K.R., Schorr, A., Johnstone, T.: *Appraisal Processes in Emotion: Theory, Methods,*
 1612 *Research.* Oxford University Press, Oxford/New York (2001)
- 1613 87. Schröder, M.: *Speech and emotion research: an overview of research frameworks and a di-*
 1614 *mensional approach to emotional speech synthesis.* Ph.D. dissertation, Univ. of Saarland,
 1615 Germany (2003)
- 1616 88. Schröder, M., Heylen, D., Poggi, I.: Perception of non-verbal emotional listener feedback.
 1617 In: Hoffmann, R., Mixdorff, H. (eds.) *Speech Prosody*, pp. 1–4 (2006)
- 1618 89. Schröder, M., Bevacqua, E., Eyben, F., Gunes, H., Heylen, D., Maat, M., Pammi, S., Pantic,
 1619 M., Pelachaud, C., Schuller, B., Sevin, E., Valstar, M., Wöllmer, M.: A demonstration of
 1620 audiovisual sensitive artificial listeners. In: *Proc. of Int. Conf. on Affective Computing and*
 1621 *Intelligent Interaction*, vol. 1, pp. 263–264 (2009)
- 1622 90. Schröder, M., Pammi, S., Gunes, H., Pantic, M., Valstar, M., Cowie, R., McKeown, G.,
 1623 Heylen, D., ter Maat, M., Eyben, F., Schuller, B., Wöllmer, M., Bevacqua, E., Pelachaud,
 1624 C., de Sevin, E.: Come and have an emotional workout with sensitive artificial listeners! In:
 1625 *Proc. of IEEE Int. Conf. on Automatic Face and Gesture Recognition* (2011)
- 1626 91. Schuller, B., Müller, R., Eyben, F., Gast, J., Hörnler, B., Wöllmer, M., Rigoll, G., Höthker,
 1627 A., Konosu, H.: Being bored? Recognising natural interest by extensive audiovisual integra-
 1628 tion for real-life application. *Image Vis. Comput.* **27**, 1760–1774 (2009)
- 1629 92. Schuller, B., Vlasenko, B., Eyben, F., Rigoll, G., Wendemuth, A.: Acoustic emotion recogni-
 1630 tion: A benchmark comparison of performances. In: *Proc. of Automatic Speech Recognition*
 1631 *and Understanding Workshop*, pp. 552–557 (2009)
- 1632 93. Schuller, B., Steidl, S., Batliner, A., Burkhardt, F., Devillers, L., Müller, C., Narayanan, S.:
 1633 The INTERSPEECH 2010 paralinguistic challenge. In: *Proc. INTERSPEECH*, pp. 2794–
 1634 2797 (2010)
- 1635 94. Schuster, M., Paliwal, K.K.: Bidirectional recurrent neural networks. *IEEE Trans. Signal*
 1636 *Process.* **45**, 2673–2681 (1997)
- 1637 95. Shen, X., Fu, X., Xuan, Y.: Do different emotional valences have same effects on spatial
 1638 attention. In: *Proc. of Int. Conf. on Natural Computation*, vol. 4, pp. 1989–1993 (2010)
- 1639 96. Sneddon, I., McKeown, G., McRorie, M., Vukicevic, T.: Cross-cultural patterns in dynamic
 1640 ratings of positive and negative natural emotional behaviour. *PLoS ONE* **6**, e14679–e14679
 1641 (2011)
- 1642 97. Soleymani, M., Davis, J., Pun, T.: A collaborative personalized affective video retrieval sys-
 1643 tem. In: *Proc. Int. Conf. on Affective Computing and Intelligent Interaction and Workshops*,
 1644 pp. 1–2 (2009)
- 1645 98. Sun, K., Yu, J., Huang, Y., Hu, X.: An improved valence-arousal emotion space for video
 1646 affective content representation and recognition. In: *Proc. IEEE Int. Conf. on Multimedia*
 1647 *and Expo*, pp. 566–569 (2009)
- 1648 99. Trouvain, J., Barry, W.J.: The prosody of excitement in horse race commentaries. In: *Proc.*
 1649 *ISCA Workshop Speech Emotion*, pp. 86–91 (2000)
- 1650 100. Truong, K.P., van Leeuwen, D.A., Neerinx, M.A. de Jong, F.M.G.: Arousal and valence
 1651 prediction in spontaneous emotional speech: Felt versus perceived emotion. In: *Proc. IN-*
 1652 *TERSPEECH*, pp. 2027–2030 (2009)
- 1653 101. Tsai, T.-C., Chen, J.-J., Lo, W.-C.: Design and implementation of mobile personal emotion
 1654 monitoring system. In: *Proc. Int. Conf. on Mobile Data Management: Systems, Services and*
 1655 *Middleware*, pp. 430–435 (2009)
- 1656 102. Tsiamyrztis, P., Dowdall, J., Shastri, D., Pavlidis, I.T., Frank, M.G., Ekman, P.: Imaging
 facial physiology for the detection of deceit. *Int. J. Comput. Vis.* (2007)
103. Wang, P., Ji, Q.: Performance modeling and prediction of face recognition systems. In: *Proc.*
of IEEE Conference on Computer Vision and Pattern Recognition, vol. 2, pp. 1566–1573
 (2006)
104. Wassermann, K.C., Eng, K., Verschure, P.F.M.J.: Live soundscape composition based on
 synthetic emotions. *IEEE Multimed.* **10**, 82–90 (2003)
105. Wöllmer, M., Eyben, F., Reiter, S., Schuller, B., Cox, C., Douglas-Cowie, E., Cowie, R.:
 Abandoning emotion classes—towards continuous emotion recognition with modelling of
 long-range dependencies. In: *Proc. INTERSPEECH*, pp. 597–600 (2008)

- 1657 106. Wöllmer, M., Schuller, B., Eyben, F., Rigoll, G.: Combining long short-term memory and
1658 dynamic Bayesian networks for incremental emotion-sensitive artificial listening. *IEEE J.*
1659 *Sel. Top. Signal Process.* **4**(5), 867–881 (2010)
- 1660 107. Yang, Y.-H., Lin, Y.-C., Su, Y.-F., Chen, H.H.: Music emotion classification: A regression
1661 approach. In: *Proc. of IEEE Int. Conf. on Multimedia and Expo*, pp. 208–211 (2007)
- 1662 108. Yu, C., Aoki, P.M., Woodruff, A.: Detecting user engagement in everyday conversations. In:
1663 *Proc. of 8th Int. Conf. on Spoken Language Processing* (2004)
- 1664 109. Zeng, Z., Pantic, M., Roisman, G.I., Huang, T.S.: A survey of affect recognition methods:
1665 Audio, visual, and spontaneous expressions. *IEEE Trans. Pattern Anal. Mach. Intell.* **31**, 39–
1666 58 (2009)
- 1667
1668
1669
1670
1671
1672
1673
1674
1675
1676
1677
1678
1679
1680
1681
1682
1683
1684
1685
1686
1687
1688
1689
1690
1691
1692
1693
1694
1695
1696
1697
1698
1699
1700
1701
1702

UNCORRECTED PROOF