



COPYRIGHT NOTICE



© 2001 IEEE. Personal use of this material is permitted. However, permission to reprint/republish this material for advertising or promotional purposes or for creating new collective works for resale or redistribution to servers or lists, or to reuse any copyrighted component of this work in other works must be obtained from the IEEE.

This material is presented to ensure timely dissemination of scholarly and technical work. Copyright and all rights therein are retained by authors or by other copyright holders. All persons copying this information are expected to adhere to the terms and constraints invoked by each author's copyright. In most cases, these works may not be reposted without the explicit permission of the copyright holder.

Toward Machine Emotional Intelligence: Analysis of Affective Physiological State

Rosalind W. Picard, *Senior Member, IEEE*, Elias Vyzas, and Jennifer Healey

Abstract—The ability to recognize emotion is one of the hallmarks of emotional intelligence, an aspect of human intelligence that has been argued to be even more important than mathematical and verbal intelligences. This paper proposes that machine intelligence needs to include emotional intelligence and demonstrates results toward this goal: developing a machine's ability to recognize human affective state given four physiological signals. We describe difficult issues unique to obtaining reliable affective data and collect a large set of data from a subject trying to elicit and experience each of eight emotional states, daily, over multiple weeks. This paper presents and compares multiple algorithms for feature-based recognition of emotional state from this data. We analyze four physiological signals that exhibit problematic day-to-day variations: The features of different emotions on the same day tend to cluster more tightly than do the features of the same emotion on different days. To handle the daily variations, we propose new features and algorithms and compare their performance. We find that the technique of seeding a Fisher Projection with the results of Sequential Floating Forward Search improves the performance of the Fisher Projection and provides the highest recognition rates reported to date for classification of affect from physiology: 81 percent recognition accuracy on eight classes of emotion, including neutral.

Index Terms—Emotion recognition, physiological patterns, feature selection, Fisher Projection, affective computing, emotional intelligence.

1 INTRODUCTION

IT is easy to think of emotion as a luxury, something that is unnecessary for basic intelligent functioning and difficult to encode in a computer program; therefore, why bother giving emotional abilities to machines? Recently, a constellation of findings, from neuroscience, psychology, and cognitive science, suggests that emotion plays surprising critical roles in rational and intelligent behavior. Most people already know that too much emotion is bad for rational thinking; much less well-known is that neuroscience studies of patients who essentially have their emotions disconnected reveal that those patients have strong impairments in intelligent day-to-day functioning, suggesting that too little emotion can impair rational thinking and behavior [1]. Apparently, emotion interacts with thinking in ways that are nonobvious but important for intelligent functioning. Emotion-processing brain regions have also been found to perform pattern recognition before the incoming signals arrive at the cortex: A rat can be taught to fear a tone even when its auditory cortex is removed, and similar emotion-oriented processing is believed to take place in human vision and audition [2].

Scientists have amassed evidence that emotional skills are a basic component of intelligence, especially for learning preferences and adapting to what is important [3], [4]. With

increasing deployment of adaptive computer systems, e.g., software agents and video retrieval systems that learn from users, the ability to sense and respond appropriately to user affective feedback is of growing importance. Emotional intelligence consists of the ability to recognize, express, and have emotions, coupled with the ability to regulate these emotions, harness them for constructive purposes, and skillfully handle the emotions of others. The skills of emotional intelligence have been argued to be a better predictor than IQ for measuring aspects of success in life [4].

Machines may never need all of the emotional skills that people need; however, there is evidence that machines will require at least some of these skills to appear intelligent when interacting with people. A relevant theory is that of Reeves and Nass at Stanford: Human-computer interaction is inherently natural and social, following the basics of human-human interaction [5]. For example, if a piece of technology talks to you but never listens to you, then it is likely to annoy you, analogous to the situation where a human talks to you but never listens to you. Nass and Reeves have conducted dozens of experiments of classical human-human interaction, taking out one of the humans and putting in a computer, and finding that the basic human-human results still hold.

Recognizing affective feedback is important for intelligent human-computer interaction. Consider a machine learning algorithm that has to decide when to interrupt the user. A human learns this by watching how you respond when you are interrupted in different situations: Did you receive the interruption showing a neutral, positive, or negative response? Without such regard for your response, the human may be seen as disrespectful, irritating, and unintelligent. One can predict a similar response toward computers that interrupt users oblivious to

• R.W. Picard is with the MIT Media Lab, 20 Ames Street, Cambridge, MA 02139. E-mail: picard@media.mit.edu.

• E. Vyzas is located at 12 Ermou Road, Vouliagmeni, Athens, Greece 16671. E-mail: evyzas@media.mit.edu.

• J. Healey is with the IBM T.J. Watson Research Center, PO Box 218, Yorktown Heights, NY 10598. E-mail: jhealey@us.ibm.com.

Manuscript received 29 Dec. 1999; revised 3 Aug. 2000; accepted 5 Mar. 2001.

Recommended for acceptance by K.W. Bowyer.

For information on obtaining reprints of this article, please send e-mail to: tpami@computer.org, and reference IEEECS Log Number 111150.

their positive or negative expressions. The computer could potentially appear more intelligent by recognizing and appropriately adapting to the user's emotion response.

Although not all computers will need emotional skills, those that interact with and adapt to humans in real-time are likely to be perceived as more intelligent if given such skills. The rest of this paper focuses on giving a machine one of the key skills of emotional intelligence: the ability to recognize emotional information expressed by a person.

1.1 Human Emotion Recognition

Human newborns show signs of recognizing affective expressions such as approval and disapproval long before they acquire language. Affect recognition is believed to play an important role in learning and in developing the ability to attend to what is important and is likely a key part of the difference between normal child development and development of autistic children, who typically have impaired affect recognition [6]. For example, instead of attending to the parent's speech with exaggerated inflection, the autistic child might tune in to an unimportant sound, missing the guidance provided by the affective cues.

Emotion modulates almost all modes of human communication—word choice, tone of voice, facial expression, gestural behaviors, posture, skin temperature and clamminess, respiration, muscle tension, and more. Emotions can significantly change the message: sometimes it is not *what* was said that was most important, but *how* it was said. Faces tend to be the most visible form of emotion communication, but they are also the most easily controlled in response to different social situations when compared to the voice and other modes of expression. Affect recognition is most likely to be accurate when it combines 1) multiple kinds of signals from the user with 2) information about the user's context, situation, goals, and preferences. A combination of low-level pattern recognition, high-level reasoning, and natural language processing is likely to provide the best emotion inference [7].

How well will a computer have to recognize human emotional state to appear intelligent? Note that no human can perfectly recognize your innermost emotions, and sometimes people cannot even recognize their own emotions. No known mode of affect communication is lossless; some aspects of internal feelings remain private, especially if you wish them to be that way or if you sufficiently disguise them. What is available to an external recognizer is what can be observed and reasoned about, and this always comes with some uncertainty. Nonetheless, people recognize each other's emotions well-enough to communicate useful feedback. Our aim is to give computers recognition abilities similar to those that people have.

1.2 Importance of Physiological Emotion Recognition

When designing intelligent machine interfaces, why not focus on facial and vocal communication—aren't these the modes that people rely upon? There are cases where such modes will be preferable, as well as other behavior-based modes, such as gestural activity or time to complete a task. However, it is a mistake to think of physiology as something that people do not naturally recognize. A stranger

shaking your hand can feel its clamminess (related to skin conductivity); a friend leaning next to you may sense your heart pounding; students can hear changes in a professor's respiration that give clues to stress; ultimately, it is muscle tension in the face that gives rise to facial expressions. People read many physiological signals of emotion.

Physiological pattern recognition of emotion has important applications in medicine, entertainment, and human-computer interaction. Affective states of depression, anxiety, and chronic anger have been shown to impede the work of the immune system, making people more vulnerable to viral infections, and slowing healing from surgery or disease ([4, chapter 11]). Physiological pattern recognition can potentially aid in assessing and quantifying stress, anger, and other emotions that influence health. Certain human physiological patterns show characteristic responses to music and other forms of entertainment, e.g., skin conductivity tends to climb as a piece of music "peps you up" and fall as it "calms you down." This principle was utilized in a wearable "affective DJ" to allow more personalized music selections than the one-size-fits-all approach of a disc jockey [8]. Changes in physiological signals can also be examined for signs of stress arising *while users interact* with technology, helping detect where the product causes unnecessary irritation or frustration, without having to interrupt the user or record her appearance. This is a new area for pattern recognition research: detecting when products cause user stress or aggravation, thereby helping developers target areas for redesign and improvement.

Physiological sensing is sometimes considered invasive because it involves physical contact with the person. However, not only is technology improving with conductive rubber and fabric electrodes that are wearable, washable, and able to be incorporated in clothes and accessories [9] but also there are new forms of noncontact physiological sensing on the horizon. In some cases, physiological sensors are perceived as less invasive than alternatives, such as video. Video almost always communicates identity, appearance, and behavior, on top of emotional information. Students engaged in distance learning may wish to communicate to the lecturer that they are furrowing their brow in confusion or puzzlement but not have the lecturer know their identity. They might not object to having a small amount of muscle tension anonymously transmitted, whereas they may object to having their appearance communicated.

New wearable computers facilitate different forms of sensing than traditional computers. Wearables often afford natural contact with the surface of the skin; however, they do not easily afford having a camera pointed at the user's face. (It can be done with a snugly fitted hat and stiff brim to mount the camera for viewing the wearer's face, a form factor that can be awkward both physically and socially.) In wearable systems, physiological sensing may be set up so that it involves no visible or heavy awkward supporting mechanisms; in this case, the physiological sensors may be less cumbersome than video sensors.

1.3 Related Research

The affect recognition problem is a hard one when you look at the few benchmarks which exist. In general, people can

recognize an emotional expression in neutral-content speech with about 60 percent accuracy, choosing from among about six different affective labels [10]. Computer algorithms match or slightly beat this accuracy, e.g., [11], [12]. Note that computer speech recognition that works at about 90 percent accuracy on neutrally-spoken speech tends to drop to 50-60 percent accuracy on emotional speech [13]. Improved handling of emotion in speech is important for recognizing what is said, as well as how it was said.

Facial expression recognition is easier for people, e.g., 70-98 percent accurate on six categories of facial expressions exhibited by actors [14] and the rates computers obtain range from 80-98 percent accuracy when recognizing 5-7 classes of emotional expression on groups of 8-32 people [15], [16]. Other research has focused not so much on recognizing a few categories of emotional expressions but on recognizing specific facial actions—the fundamental muscle movements that comprise Paul Ekman's Facial Action Coding System—which can be combined to describe all facial expressions. Recognizers have already been built for a handful of the facial actions [17], [18], [19], [20], and the automated recognizers have been shown to perform comparably to humans trained in recognizing facial actions [18]. These facial actions are essentially *facial phonemes*, which can be assembled to form facial expressions. There are also recent efforts that indicate that combining audio and video signals for emotion recognition can give improved results [21], [22], [23].

Although the progress in facial, vocal, and combined facial/vocal expression recognition is promising, all of the results above are on presegmented data of a small set of sometimes exaggerated expressions or on a small subset of hand-marked singly-occurring facial actions. The state-of-the-art in affect recognition is similar to that of speech recognition several decades ago when the computer could classify the carefully articulated digits, "0, 1, 2, . . . , 9," spoken with pauses in between, but could not accurately detect these digits in the many ways they are spoken in larger continuous conversations.

Emotion recognition research is also hard because understanding emotion is hard; after over a century of research, emotion theorists still do not agree upon what emotions are and how they are communicated. One of the big questions in emotion theory is whether distinct physiological patterns accompany each emotion [24]. The physiological muscle movements comprising what looks to an outsider to be a facial expression may not always correspond to a real underlying emotional state. Emotion consists of more than its outward physical expression; it also consists of internal feelings and thoughts, as well as other internal processes of which the person having the emotion may not be aware.

The relation between internal bodily feelings and externally observable expression is still an open research area, with a history of controversy. Historically, James was the major proponent of emotion as an experience of bodily changes, such as your heart pounding or your hands perspiring [25]. This view was challenged by Cannon [26] and again by Schachter and Singer who argued that the experience of physiological changes was not sufficient to

discriminate emotions. Schachter and Singer's experiments showed that, if a bodily arousal state was induced, then subjects could be put into two distinct moods simply by being put in two different situations. They argued that physiological responses such as sweaty palms and a rapid heart beat inform our brain that we are aroused and then the brain must appraise the situation we are in before it can label the state with an emotion such as fear or love [27].

Since the classic work of Schachter and Singer, there has been a debate about whether or not emotions are accompanied by specific physiological changes other than simply arousal level. Ekman et al. [28] and Winton et al. [29] provided some of the first findings showing significant differences in autonomic nervous system signals according to a small number of emotional categories or dimensions, but there was no exploration of automated classification. Fridlund and Izard [30] appear to have been the first to apply pattern recognition (linear discriminants) to classification of emotion from physiological features, attaining rates of 38-51 percent accuracy (via cross-validation) on subject-dependent classification of four different facial expressions (happy, sad, anger, fear) given four facial electromyogram signals. Although there are over a dozen published efforts aimed at finding physiological correlates when examining small sets of emotions (from 2-7 emotions according to a recent overview [31]), most have focused on t-test or analysis of variance comparisons, combining data over many subjects, where each was measured for a relatively small amount of time (seconds or minutes). Relatively few of the studies have included neutral control states where the subject relaxed and passed time feeling no specific emotion, and none to our knowledge have collected data from a person repeatedly, over many weeks, where disparate sources of noise enter the data. Few efforts beyond Fridlund's have employed linear discriminants, and we know of none that have applied more sophisticated pattern recognition to physiological features.

The work in this paper is novel in trying to classify physiological patterns for a set of eight emotions (including neutral), by applying pattern recognition techniques beyond that of simple discriminants to the problem (we use new features, feature selection, spatial transformations of features, and combinations of these methods) and by focusing on "felt" emotions of a single subject gathered over sessions spanning many weeks. The results we obtain are also independent of psychological debates on the universality of emotion categories [32], focusing instead on user-defined emotion categories.

The contributions of this paper include not only a new means for pattern analysis of affective states from physiology, but also the finding of significant classification rates from physiological patterns corresponding to eight affective states measured from a subject over many weeks of data. Our results also reveal significant discrimination among both most commonly described dimensions of emotion: valence and arousal. We show that the day-to-day variations in physiological signals are large, even when the same emotion is expressed, and this effect undermines recognition accuracy if it is not appropriately handled. This paper proposes and compares techniques for handling

day-to-day variation and presents new results in affect recognition based on physiology. The results lie between the rates obtained for expression recognition from vocal and facial features and are the highest reported to date for classifying eight emotional states given physiological patterns.

2 GATHERING GOOD AFFECTIVE DATA

In computer vision or speech recognition, it has become easy to gather meaningful data; frame-grabbers, microphones, cameras, and digitizers are reliable, easy to use, and the integrity of the data can be seen or heard by non-specialists; however, nonspecialists do not usually know what comprises a good physiological signal. Although people recognize emotional information from physiology, it is not natural to do so by looking at 1D signal waveforms. Not only does it take effort to learn what a good signal is, but the sensing systems (A/D converters and buffering-data capture systems) for physiology do not seem to be as reliable as those for video and audio. Factors such as whether or not the subject just washed her hands, how much gel she applied under an electrode, motion artifacts, and precisely where the sensor was placed, all affect the readings. These are some of the technical factors that contribute to the difficulty in gathering accurate physiological data.

Although dealing with the recording devices can be tricky, a much harder problem is that of obtaining the ground truth of the data, or getting data that genuinely corresponds to a particular emotional state. In vision or speech research, the subject matter is often objective: scene depth, words spoken, etc. In those cases, the ground-truth labels for the data are easily obtained. Easy-to-label data is sometimes obtained in emotion research when a singular strong emotion is captured, such as an episode of rage. However, more often, the ground truth—which emotion was present—is difficult to establish.

Consider a task where a person uses the computer to retrieve images. Suppose our job is to analyze physiological signals of the user as he or she encounters pleasing or irritating features or content within the system. We would like to label the data according to the emotional state of the user (ground truth). Here, the problem is complicated because there is little way of knowing whether the person was truly pleased or irritated when encountering the stimuli intended to induce these states. We may have tried to please, but the user was irritated because of something she remembered unrelated to the task. We may have tried to irritate and not succeeded. If we *ask* the person how she felt, her answer can vary according to her awareness of her feelings, her comfort in talking about feelings, her rapport with the administrator(s) of the experiment, and more. *When you ask* is also important—soon and often is likely to be more accurate, but also more irritating, thereby changing the emotional state. Thus, measurement of ground truth disturbs the state of that truth. When it comes to faces or voices, we can see if the person was smiling or hear if her voice sounded cheerful, but that still does not mean that she was happy. With physiology, little is known about *how* emotions make their impact, but the signals are also

potentially more sincere expressions of the user's state since they tend to be less mediated by cognitive and social influences.

2.1 Five Factors in Eliciting Emotion

In beginning this research, we thought it would be simple to ask somebody to feel and express an emotion and to record data during such episodes. Indeed, this is the approach taken in most facial and vocal expression studies to date: turn the camera and microphone on, ask the subject to express joy, anger, etc., record it, and label it by what was requested. However, obtaining high quality physiological data for affect analysis requires attention to experimental design issues not traditionally required in pattern recognition research. Here, we outline five (not necessarily independent) factors that influence data collection, to serve as a useful guide for researchers trying to obtain affect data. We summarize the factors by listing their extreme conditions, but there are also in-between conditions:

1. *Subject-elicited* versus *event-elicited*: Does subject purposefully elicit emotion or is it elicited by a stimulus or situation outside the subject's efforts?
2. *Lab setting* versus *real-world*: Is subject in a lab or in a special room that is not their usual environment?
3. *Expression* versus *feeling*: Is the emphasis on external expression or on internal feeling?
4. *Open-recording* versus *hidden-recording*: Does subject know that anything is being recorded?¹
5. *Emotion-purpose* versus *other-purpose*: Does subject know that the experiment is about emotion?

The most natural setup for gathering genuine emotions is opportunistic: The subject's emotion occurs as a consequence of personally significant circumstances (*event-elicited*); it occurs while they are in some natural location for them (*real-world*); the subject feels the emotion internally (*feeling*); subject behavior, including expression, is not influenced by knowledge of being in an experiment or being recorded (*hidden-recording, other-purpose*). Such data sets are usually impossible to get because of privacy and ethics concerns, but as recording devices are increasingly prevalent, people may cease to be aware of them, and the data captured by these devices can be as if the devices were hidden. Researchers may try to create such opportunistic situations; for example, showing an emotion-eliciting film in a theater without telling subjects the true purpose of the showing, and without telling them that they have been videotaped until afterward. Even so, an emotion-inducing

1. When the presence of the camera or other recording device is hidden, it is ethically necessary (via the guidelines of the MIT Committee on the Use of Humans as Experimental Subjects) to debrief the subject, let them know why the secrecy was necessary, tell them what signals have been recorded, and obtain their permission for data analysis, destroying the data if the subject withholds permission. Unlike with video, it is not yet possible to record a collection of physiological signals without the subject being aware of the sensors in some form (although this is changing with new technology.) However, subjects can be led to believe that their physiology is being sensed for some reason other than emotions, which is often an important deception since a subject who knows you are trying to make them frustrated may, therefore, not get frustrated. We have found that such deception is acceptable to almost all subjects when properly conducted. Nonetheless, we believe deception should not be used unless it is necessary and then only in accord with ethical guidelines.

movie scene will affect some viewers more than others and some maybe not at all.

The opportunistic situation contrasts with the experiments typically used to gather data for facial or vocal expression recognition systems. The common setup is one in which data of a *subject-elicited external expression* in a *lab setting*, in front of a visible *open-recording* camera or microphone, and with knowledge that the data will be used for analysis of emotion (*emotion-purpose*) are gathered. Such expressions, made with or without corresponding internal feelings, are “real” and important to recognize, even if not accompanied by true feelings. Social communication involves both unfelt emotions (emphasizing expression) and more genuine ones (emphasizing feeling.) A protocol for gathering data may emphasize external expression or internal feeling; both are important. Moreover, there is considerable overlap since physical expression can help induce the internal feeling, and vice-versa.

In this paper, we gathered real data following a *subject-elicited*, close to *real-world* (subject’s comfortable usual workplace), *feeling*, *open-recording*, and *emotion-purpose* methodology. The key one of these factors that makes our data unique is that the subject tried to elicit an internal *feeling* of each emotion.

2.2 Single-Subject Multiple-Day Data Collection

The data we gather is from a single subject over many weeks of time, standing in contrast to efforts that examine many subjects over a short recording interval (usually single session on only one day). Although the scope is limited to one subject, the amount of data for this subject encompasses a larger set than has traditionally been used in affect recognition studies involving multiple subjects. The data are potentially useful for many kinds of analysis and will be made available for research purposes.

There are many reasons to focus on person-dependent recognition in the early stages of affect recognition, even though some forms of emotion communication are not only person-independent, but have been argued, namely, by Paul Ekman, to be basic across different cultures. Ekman and colleagues acknowledge that even simply labeled emotions like “joy” and “anger” can have different interpretations across individuals within the same culture; this complicates the quest to see if subjects elicit similar physiological patterns for the same emotion. When lots of subjects have been examined over a short amount of time, researchers have had difficulty finding significant physiological patterns, which may be in part because physiology can vary subtly with how each individual interprets each emotion. By using one subject, who tried to focus on the same personal interpretation of the emotion in each session, we hoped to maximize the chance of getting consistent interpretations for each emotion. This also means that the expressive data can be expected to differ for another subject since the way another subject interprets and reacts to the emotions may differ. Hence, a weakness of this approach is that the precise features and recognition results we obtain with this data may not be the same for other subjects. However, the methodology for gathering and analyzing the data in this paper is not dependent on the subject; the approach described in this paper is general.

Emotion recognition is in an early stage of research, similar to early research on speech recognition, where it is valuable to develop person-dependent methods. For personal computing applications, we desire the machine to learn an individual’s patterns and not just some average response formed across a group, which may not apply to the individual.

The subject in our experiments was a healthy graduate student with two years acting experience plus training in visualization, who was willing to devote over six weeks to data collection. The subject sat in her quiet workspace early each day, at roughly the same time of day, and tried to experience eight affective states with the aid of a computer controlled prompting system, the “Sentograph,” developed by Clynes [33] and a set of personally-significant imagery she developed to help elicit the emotional state.

The Clynes protocol for eliciting emotion has three features that contribute to helping the subject feel the emotions and that make it appropriate for physiological data collection: 1) It sequences eight emotions in a way that supposedly makes it easier for many people to transition from emotion to emotion. 2) It engages physical expression—asking the subject to push a finger against a button with a dual axis pressure sensor in an expressive way—but in a way that limits motion artifacts being introduced to the physiological signals. Physical expression gives somatosensory feedback to the subject, a process that can help focus and strengthen the feeling of the emotion [34]. 3) It prompts the subject to repeatedly express the same emotion during an approximately three minute interval, at a rate dependent on the emotion in order to try to intensify the emotional experience [33].

The order of expression of the eight states: *no emotion*, *anger*, *hate*, *grief*, *platonic love*, *romantic love*, *joy*, and *reverence* was found by Clynes to help subjects reliably feel each emotion; for example, it would probably be harder on most subjects to have to oscillate between the positive and negative states, e.g., *joy*, *hate*, *platonic love*, *grief*, etc. However, because interpretations for each of the emotions may vary with each individual, we do not expect this order to be optimal for everyone.

Descriptive guidelines on the meaning of each emotion were developed by the subject before the experiment. The subject reported the images she used to induce each state, the degree to which she found each experience arousing (exciting, distressing, disturbing, tranquil), and the degree to which she felt the emotion was positive or negative (valence) (See Table 1). Daily ratings varied in intensity—both up and down, with no particular trend, but the overall character of each state was consistent over the weeks of data collection.

The eight emotions in this study differ from those typically explored. Although there is no widespread agreement on the definition and existence of “basic” emotions and which names would comprise such a list, researchers on facial expressions tend to focus on *anger*, *fear*, *joy*, and *sadness* with *disgust* and *surprise* often examined as well as *contempt*, *acceptance*, and *anticipation* (e.g., [32]). Theorists still do not agree on what an emotion is and many of them do not consider *love* and *surprise* to be emotions.

TABLE 1
The Subject's Descriptions of Imagery and Emotional Character Used for Each of the Eight Emotions

Emotion	Imagery	Description	Arousal	Valence
(N)o emotion	blank paper, typewriter	boredom, vacancy	low	neut.
(A)nger	people who arouse rage	desire to fight	very high	very neg.
(H)ate	injustice, cruelty	passive anger	low	neg.
(G)rief	deformed child, loss of mother	loss, sadness	high	neg.
(P)latoic love	family, summer	happiness, peace	low	pos.
Romantic (L)ove	romantic encounters	excitement, lust	very high	pos.
(J)oy	The music "Ode to Joy"	uplifting happiness	med. high	pos.
(R)everence	church, prayer	calm, peace	very low	neut.

Our work is less concerned with finding a set of "basic" emotions and more concerned with giving computers the ability to recognize whatever affective states might be relevant in a personalized human-computer interaction. The ideal states for a computer to recognize will depend on the application. For example, in a learning-tutor application, detecting expressions of curiosity, boredom, and frustration may be more relevant than detecting emotions on the theorists' "basic" lists.

Clynes' set of eight was motivated by considering emotions that have been communicated through centuries of musical performance on several continents. We started with his set not because we think it is the best for computer-human interaction (such a set is likely to vary with computer applications—entertainment, business, socializing, etc.), but rather because this set together with its method for elicitation had shown an ability to help subjects reliably feel the emotions and had shown repeatable signs of physical differentiation in how subjects' finger pressure applied to a finger rest differs with each emotion [33], [35], a measurable outcome that suggests different states were being achieved. It was important to our investment in long-term data collection that we have a reliable method of helping the user repeatedly generate distinct emotional states.

For the purposes of this research, the specific emotions and their definitions are not as important as the fact that 1) the subject could relate to the named emotion in a consistent, specific, and personal way and 2) the emotion categories span a range of high and low *arousal* and positive and negative *valence*. These two dimensions are believed to be the most important dimensions for categorizing emotions [36] and continue to be used for describing emotions that arise in many contexts, including recent efforts to categorize emotions arising when people look at imagery [37]. The arousal axis ranges from calm and peaceful to active and excited, while the valence axis ranges from negative (displeasing) to positive (pleasing).

2.3 Experimental Method and Construction of Data Sets

Data were gathered from four sensors: a triode electromyogram (\mathcal{E}) measuring facial muscle tension along the masseter (with Ag-AgCl electrodes of size 11mm each and 10-20 high-conductivity gel), a photoplethysmograph measuring blood volume pressure (\mathcal{B}) placed on the tip of

the ring finger of the left hand, a skin conductance (\mathcal{S}) sensor measuring electrodermal activity from the middle of the three segments of the index and middle fingers on the palm-side of the left hand (with 11mm Ag-AgCl electrodes and K-Y Jelly used for low-conductivity gel), and a Hall effect respiration sensor (\mathcal{R}) placed around the diaphragm. The left hand was held still throughout data collection and the subject was seated and relatively motionless except for small pressure changes she applied with her right hand to the finger rest. Sensors and sampling were provided by the Thought Technologies ProComp unit, chosen because the unit is small enough to attach to a wearable computer and offers eight optically isolated channels for recording. Signals were sampled at 20 Hz.² The ProComp automatically computed the heart rate (\mathcal{H}) as a function of the inter-beat intervals of the blood volume pressure, \mathcal{B} . More details on this system and on our methodology are available [38].

Each day's session lasted around 25 minutes, resulting in around 28 to 33 thousand samples per physiological signal, with each different emotion segment being around two to five thousand samples long, due to the variation built into the Clynes method of eliciting the emotional states [33]. Eight signal segments of the raw data (2,000 samples each) from Data Set I are shown in Fig. 1. On roughly a third of the 30 days for which we collected data, either one or more sensors failed during some portion of the 25-minute experiment because an electrode came loose or one or more channels failed to sample and save some of the data properly. From the complete or nearly-complete sessions, we constructed two overlapping Data Sets.

Data Set I was assembled before the 30 days were over, and was formed as follows: Data segments of 2,000 samples (100 seconds) in length were taken from each of the signals \mathcal{E} , \mathcal{B} , \mathcal{G} , and \mathcal{R} for each of the eight emotions, on each of 19 days where there were no failures in these segments of data collection. The 2,000 samples were taken from the end of each emotion segment to avoid the transitional onset where the subject was prompted to move to the next emotion. A 20th day's data set was created out of a combination of partial records in which some of the sensors had failed.

2. The electromyogram is the only signal for which this sampling rate should have caused aliasing. However, our investigation of the signal showed that it registered a clear response when the jaw was clenched versus relaxed; thus, it was satisfactory for gathering coarse muscle tension information.

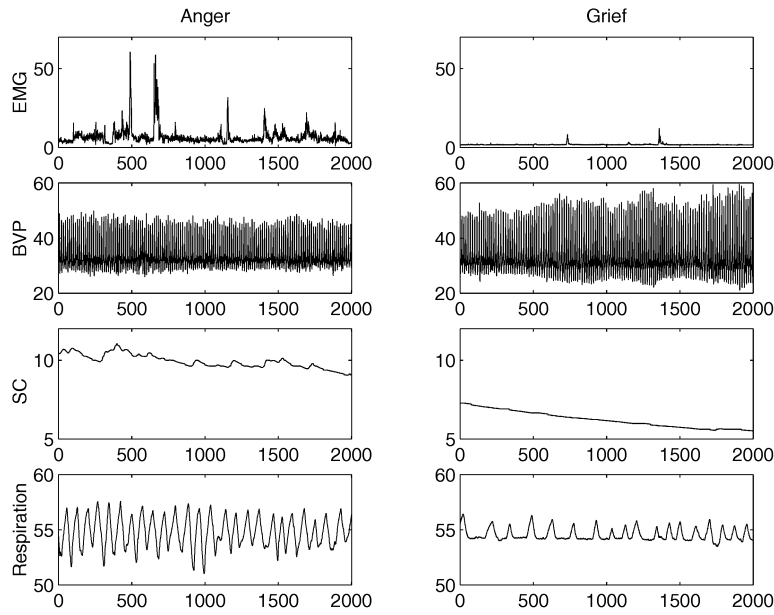


Fig. 1. Examples of physiological signals measured from a user while she intentionally expressed anger (left) and grief (right). From top to bottom: electromyogram (microvolts), blood volume pressure (percent reflectance), skin conductivity (microSiemens), and respiration (percent maximum expansion). Each box shows 100 seconds of response. The segments shown here are visibly different for the two emotions, which was not true in general.

Data Set II is a larger data set comprised of 20 days in which the sensors did not fail during any part of the experiment for the five signals: \mathcal{E} , \mathcal{B} , \mathcal{G} , \mathcal{R} , and \mathcal{H} . These 20 days included 16 of the original days from Data Set I, but, for all 20 days, we used all of the samples available for each emotion, thereby including the transitional regions. Because different emotions lasted for different lengths of time, the 2,000 samples in Data Set I were at times closer or farther away from the beginning of an emotion segment. To avoid this bias and to maximize data available for training and testing, Data Set II includes all the samples for all the emotions and signals over 20 days, roughly 2,000 to 5,000 samples per emotion per signal per day. With an average of twice the number of samples (or minutes of data) as Data Set I, Data Set II resulted in an average 10 percent gain in performance when we compared it to Data Set I across all the methods.

3 FEATURE EXTRACTION, SELECTION, AND TRANSFORMATION

Because the signals involved have different and complex sources, because there are not yet good models to describe them, and because of an interest in seeing how some classical methods perform before an extensive modeling effort is launched, we choose in this paper to explore a feature-based approach to classification.

3.1 Proposed Feature Sets

The psychophysiology and emotion literature contains several efforts to identify features of bodily changes (facial muscle movements, heartrate variations, etc.) that might correlate with having an emotion (e.g., [31]). We gather a variety of features, some from the literature and some that we propose. Several that we propose are physically

motivated, intended to capture the underlying nature of specific signals, such as the way respiration is quasi-periodic, while others are simple statistics and nonlinear combinations thereof. We do not expect the best classifier to require all the features proposed below or even to require such a huge number of features. Our effort is to advance the state-of-the-art in pattern recognition of affect from physiology by proposing a large space of reasonable features and systematically evaluating subsets of it and transformations thereof. Below, six statistical features are presented, followed by 10 more physically-motivated features aimed at compensating for day-to-day variations. Which features were found to be most useful (as a function of each of the classifiers) will be summarized later in Table 9.

The six statistical features can be computed for each of the signals as follows: Let the signal $\{\mathcal{E}, \mathcal{B}, \mathcal{G}, \mathcal{R}, \mathcal{H}\}$ from any one of the eight emotion segments be designated by X . The signal is gathered for eight different emotions each day, for 20 days. Let X_n represent the value of the n th sample of the raw signal, where $n = 1, \dots, N$, with $N = 2,000$ for Data Set I, and with N in the range of 2,000 to 5,000 for Data Set II. Let \tilde{X}_n refer to the normalized signal (zero mean, unit variance):

$$\tilde{X}_n = \frac{X_n - \mu_X}{\sigma_X} \quad i = 1, \dots, 4,$$

where μ_X and σ_X are the means and standard deviations of X as explained below. Following are six statistical features we investigated:

1. the means of the raw signals

$$\mu_X = \frac{1}{N} \sum_{n=1}^N X_n, \tag{1}$$

2. the standard deviations of the raw signals,

$$\sigma_X = \left(\frac{1}{N-1} \sum_{n=1}^N (X_n - \mu_X)^2 \right)^{1/2}, \quad (2)$$

3. the means of the absolute values of the first differences of the raw signals

$$\delta_X = \frac{1}{N-1} \sum_{n=1}^{N-1} |X_{n+1} - X_n|, \quad (3)$$

4. the means of the absolute values of the first differences of the normalized signals

$$\tilde{\delta}_X = \frac{1}{N-1} \sum_{n=1}^{N-1} |\tilde{X}_{n+1} - \tilde{X}_n| = \frac{\delta_X}{\sigma_X}, \quad (4)$$

5. the means of the absolute values of the second differences of the raw signals

$$\gamma_X = \frac{1}{N-2} \sum_{n=1}^{N-2} |X_{n+2} - X_n|, \quad (5)$$

6. the means of the absolute values of the second differences of the normalized signals

$$\tilde{\gamma}_X = \frac{1}{N-2} \sum_{n=1}^{N-2} |\tilde{X}_{n+2} - \tilde{X}_n| = \frac{\gamma_X}{\sigma_X}. \quad (6)$$

The features (1), (2), (3), (4), (5), and (6) were chosen to cover and extend a range of typically measured statistics in the emotion physiology literature [39]. (Means and variances are already commonly computed; the first difference approximates a gradient.) Note that not all the features are independent; in particular, $\tilde{\delta}_X$ and $\tilde{\gamma}_X$ are nonlinear combinations of other features. Also, the heart rate signal, \mathcal{H} , is derived from the blood volume pressure signal, \mathcal{B} , by a nonlinear transformation performed automatically by the ProComp sensing system. It did not require an additional sensor and, so, was dependent upon \mathcal{B} . The dependencies are not linear; consequently, they are not obtainable by the linear combination methods used later to reduce the dimensionality of the feature space and can potentially be of use in finding a good transformation for separating the classes. The comparisons below will verify this.

One advantage of the features in (1), (2), (3), (4), (5), and (6) is that they can easily be computed in an online way [40], which makes them advantageous for real-time recognition systems. However, the statistical features do not exploit knowledge we may have about the physical sources of the signals, and provide no special normalization for day-to-day variations in the signals. Factors such as hand washing, gel application, and sensor placement can easily affect the statistics. These influences combine with the subject's daily mood and with other cognitive and bodily influences in presently unknown ways, making them hard to model.

In an effort to compensate for some of the nonemotion-related variations of the signals and to include more physically-motivated knowledge about the signal (such as underlying periodic excitations), we also compute and evaluate another set of 10 physiology-dependent features, $f_1 - f_{10}$, described below.

From the interbeat intervals of the blood volume pressure waveform, the Procomp computes the heart rate, \mathcal{H} , which is approximately 1/the interbeat intervals. We applied a 500-point (25 sec) Hanning window, h , [41] to form a smoothed heartbeat rate, $b = \mathcal{H} * h$, then took the mean:

$$f_1 = \frac{1}{N} \sum_{n=1}^N b_n. \quad (7)$$

The average acceleration or deceleration of the heart beat rate was calculated by taking the mean of the first difference:

$$f_2 = \frac{1}{N-1} \sum_{n=1}^{N-1} (b_{n+1} - b_n) = \frac{1}{N-1} (b_N - b_1). \quad (8)$$

The skin conductivity signal, \mathcal{S} contains high frequency fluctuations that may be noise; these fluctuations are reduced by convolving with h , a 25 second Hanning window, to form $s = \mathcal{S} * h$. We also use a form of contrast normalization to account for baseline fluctuations; this measure was proposed by Rose [42] and found to be valuable over years of psychophysiology [43], [44], where $\max(g)$ and $\min(g)$ are, with respect to the whole day's data (for all emotions that day):

$$f_3 = \frac{s \min(s)}{\max(s) - \min(s)}. \quad (9)$$

The mean of the first difference of the smoothed skin conductivity is also proposed:

$$f_4 = \mu_{(s_{n+1} - s_n)} = \frac{1}{N-1} (s_N - s_1). \quad (10)$$

The respiration sensor measured expansion and contraction of the chest cavity using a Hall effect sensor attached around the chest with a velcro band. Let N_d be the number of samples collected that day. To account for variations in the initial tightness of the sensor placement from day to day, we formed the mean of the whole day's respiration data:

$$\mu_{\mathcal{R}, \text{day}} = \frac{1}{N_d} \sum_{n=1}^{N_d} \mathcal{R}_n, \quad (11)$$

and then subtracted this to get $r = \mathcal{R} - \mu_{\mathcal{R}, \text{day}}$. Two respiration features were then formed as:

$$f_5 = \frac{1}{N} \sum_{n=1}^N r_n \quad (12)$$

and

$$f_6 = \frac{1}{N-1} \sum_{n=1}^N (r_n - \mu_{\mathcal{R}, \text{day}})^2. \quad (13)$$

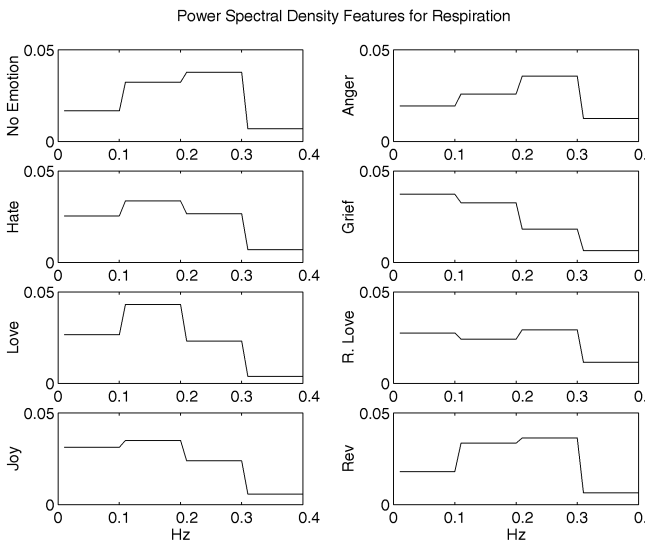


Fig. 2. Power spectral density from 0.0-0.4 Hz in the respiration signal for eight emotions. The heights of the four bins shown here were used as features $f_7 - f_{10}$.

The four features $f_7 - f_{10}$ represent frequency information in the respiration signal. Each was computed using the power spectral density function (PSD command) of Matlab, which uses Welch's averaged periodogram. The features, illustrated in Fig. 2, represent the average energy in each of the first four 0.1 Hz bands of the power spectral density range 0.0-0.4 Hz.

3.2 Selecting and Transforming Features

We first compare two techniques that have appeared in the literature to establish a benchmark: Sequential Floating Forward Search (SFFS) and Fisher Projection (FP). Next, we propose and compare a new combination of these, which we label SFFS-FP. This combination is motivated by the dual roles of the two methods: SFFS *selects* from a set of features, while Fisher Projection *linearly transforms* a set of features. Since feature selection is nonlinear, the cascade of these two methods should provide a more powerful combination than either alone. This was confirmed in our experiments.

The **Sequential Floating Forward Search (SFFS)** method [45] is chosen because of its consistent success in previous evaluations of feature selection algorithms, where it has been shown to outperform methods such as Sequential Forward Search (SFS), Sequential Backward Search (SBS), Generalized SFS and SBS, and Max-Min in comparison studies [46]. Of course the performance of SFFS is data dependent and the data set here is new and different; SFFS may not be the best method to use. Nonetheless, because of its well-documented success in other pattern recognition problems, it will help establish a benchmark for this new application area. The SFFS method takes as input the values of n features. It then does a nonexhaustive search on the feature space by iteratively including and omitting features. It outputs one subset of m features for each m , $2 \leq m \leq n$, together with its classification rate. The algorithm is described in detail in [47].

Fisher Projection (FP) [48] is a well-known method of reducing dimensionality by finding a linear projection of the data to a space of fewer dimensions where the classes are well-separated. Due to the nature of the Fisher projection method, the data can only be projected down to $c - 1$ (or fewer if one wants) dimensions, assuming that originally there are more than $c - 1$ dimensions and c is the number of classes. If the amount of training data is inadequate, or if the quality of some of the features is questionable, then some of the dimensions of the Fisher projection may be a result of noise rather than a result of differences among the classes. In this case, Fisher might find a meaningless projection which reduces the error in the training data but performs poorly in the testing data. For this reason, we not only separate training and testing data, but we also evaluate projections down to fewer than $c - 1$ dimensions. Note that if the number of features n is smaller than the number of classes c , the Fisher projection is meaningful only up to at most $n - 1$ dimensions. Therefore, the number of Fisher projection dimensions d is $1 \leq d \leq \min(n, c) - 1$, e.g., when 24 features are used on all eight classes, all $d = [1, 7]$ are tried, and when four features are used on eight classes, all $d = [1, 3]$ are tried.

A **Hybrid SFFS with Fisher Projection (SFFS-FP)** method is proposed, implemented, and evaluated here for comparison. As mentioned above, the SFFS algorithm proposes one subset of m features for each m , $2 \leq m \leq n$. It *selects*, but does not *transform* the features. Instead of feeding the Fisher algorithm with all possible features, we use the subsets that the SFFS algorithm proposes as input to the Fisher Algorithm. Note that the SFFS method is used here as a preprocessor for reducing the number of features fed into the Fisher algorithm, and not as a classification method.

4 CLASSIFICATION

This section describes and compares the results of a set of classification experiments leading up to the best results (81 percent) shown in Tables 6 and 7.

The SFFS software employs a k-nearest-neighbor (k-NN) classifier [49], so that it not only outputs the best set of features according to this classifier, but their classification accuracy as well. We used the k-NN classifier for benchmarking the SFFS method, following the methodology of Jain and Zongker [46]. For FP and SFFS-FP, we used a MAP classifier, with details below.

In all three comparisons, SFFS, FP, and SFFS-FP, we used the leave-one-out method for cross-validation because of the relatively small amount of data available and the high dimensional feature spaces. In each case, the data point (vector of features for one day's data for one emotion) to be classified was excluded from the data set before the SFFS, FP, or SFFS-FP was run. The best set of features or best transform (determined from the training set) was then applied to the test data point to determine classification accuracy.

For each k , where we varied $1 \leq k \leq 20$, the SFFS algorithm output one set of m features for each $2 \leq m \leq n$. For SFFS-FP, we computed all possible Fisher projections for each of these feature sets. For both the FP and SFFS-FP methods, we then

TABLE 2

Classification Accuracy for Three Methods Applied to Data Set I, Starting with the 24 Features $\mu_X, \sigma_X, \delta_X, \tilde{\delta}_X, \gamma_X, \tilde{\gamma}_X, X \in (\mathcal{E}, \mathcal{B}, \mathcal{G}, \mathcal{R})$

Number of Emotions	Random Guessing (%)	SFFS (%)	Fisher (%)	SFFS-FP (%)
8	12.5	40.6	40.0	46.3
5 (NAGJR)	20.0	64.0	60.0	65.0
4 (NAGR)	25.0	70.0	61.3	68.7
4 (AGJR)	25.0	72.5	60.0	67.5
3 (AGR)	33.3	83.3	71.7	80.0
3 (AJR)	33.3	88.3	66.7	83.3

The best-recognized emotions are denoted by their first initial: (N)neutral, (A)nger, (G)rief, (J)oy, (R)everence.

fit Gaussians to the data in the reduced d -dimensional feature space and applied Maximum a Posteriori (MAP) classification to the features in the Fisher-transformed space. The specific algorithm was:

1. The data point to be classified (the testing set only includes one point) is excluded from the data set. The remaining data set is used as the training set.
2. The Fisher projection matrix (used in either FP or SFFS-FP) is calculated from only the training set. Then, both the training and testing set are projected down to the d dimensions found by Fisher.
3. The data in the d -dimensional space is assumed to be Gaussian. The respective means and covariance matrices of the classes are estimated from the training data.
4. The posterior probability of the data point is calculated.
5. The data point is then classified as coming from the class with the highest posterior probability.

4.1 Initial Results—Data Set I

Our first comparison was made on Data Set I using (1), (2), (3), (4), (5), and (6) computed on the four raw physiological signals, yielding 24 features: $\mu_X, \sigma_X, \delta_X, \tilde{\delta}_X, \gamma_X, \tilde{\gamma}_X, X \in (\mathcal{E}, \mathcal{B}, \mathcal{G}, \mathcal{R})$. The results of SFFS, Fisher, and SFFS-FP are shown in Table 2. Table 2 shows the results applied to all eight emotions, as well as the results applied to all sets of C emotions, where $C = 3, 4, 5$. The best-recognized emotions were (N)neutral, (A)nger, (G)rief, (J)oy, and (R)everence. The states of (H)ate and the states of love, (P)latic Love and Romantic (L)ove, were not well-discriminated by any of the classifiers.

How good are these results—are the differences between the classifiers significant? For each pair of results here and through the rest of the paper, we computed the probability that the lower error rate really is lower, treating error rate over the 160 (or fewer) trials as a Bernoulli random variable. For example, (last row of Table 2) we compared the 20 errors made by Fisher to the 10 errors made by SFFS-FP, over the 60 trials for the three emotions (AJR). The confidence that the performance was improved was 98 percent. Although the performance improvements of SFFS-FP over Fisher in the other rows range from 5 to over 8 percentage points, the confidences range from 77 percent to 87 percent; thus, the

TABLE 3

Number of Dimensions and Number of Features m Used in the Fisher Projections and Feature Selection Algorithms that Gave the Results of Table 2

Number of Emotions	No. of Dimensions		No. of Features	
	Fisher-24	SFFS-FP	m_{SFFS}	$m_{SFFS-FP}$
8	6/7	4,5/7	13	17
5 (NAGJR)	3/4	3/4	12-17	15
4 (NAGR)	3/3	3/3	9-15,18	19
4 (AGJR)	3/3	2/3	7-8	12
3 (AGR)	2/2	2/2	2-16	12
3 (AJR)	1/2	2/2	6-14	7

The denominator in the dimensions columns is the maximum number of dimensions that could have been used; half of the time the results were better using a number less than this maximum (the numerator). When a range is shown for the number of features, m , then the performance was the same for the whole range.

performance improvement of SFFS-FP cannot be said to be statistically significant in these cases, given that 90-95 percent confidence is usually required for that statement. Nonetheless, all the classifiers in this table provide statistically significant improved performance over a random classifier (confidence > 99.99 percent).

Table 3 shows the number of dimensions and the number of features, respectively, that gave the best results. From looking at these numbers, we can conclude that performance was sometimes improved by projecting down to fewer than $c - 1$ dimensions for a c class problem. We can also see that a broad range of numbers of features led to the best results, but, in no case, were 20 or more features useful for constructing the best classifier.

4.2 The Problem of Day-Dependence

In visually examining several projections into 2D, we noticed that the features of different emotions from the same day often clustered more closely than did features for the same emotions on different days. To try to quantify this “day dependence,” we ran a side experiment: How hard would it be to classify *what day* each emotion came from, as opposed to *which emotion* was being expressed? Because classifying days is not the main interest, we only ran one comparison: Fisher Projection applied to the 24 features $\mu_X, \sigma_X, \delta_X, \tilde{\delta}_X, \gamma_X, \tilde{\gamma}_X, X \in (\mathcal{E}, \mathcal{B}, \mathcal{G}, \mathcal{R})$ computed on Data Set I. Leave-one-out cross-validation was applied to every point, and MAP classification was used on the classes (days or emotions), as described above. When the classes were $c = 20$ days (versus the $c = 8$ emotions), then the recognition accuracy jumped to 83 percent, which is significantly better than random guessing (5 percent) and significantly better than the 40 percent found for emotion classification using these features.

The day dependence is likely due to three factors: 1) skin-sensor interface influences—including hand washing, application of slightly different amounts of gel, and slight changes in positioning of the sensors; 2) variations in physiology that may be caused by caffeine, sugar, sleep, hormones, and other nonemotional factors; 3) variations in physiology that are mood and emotion dependent—such as an inability to build up an intense experience of joy if the subject felt a strong baseline mood of sadness that day. The

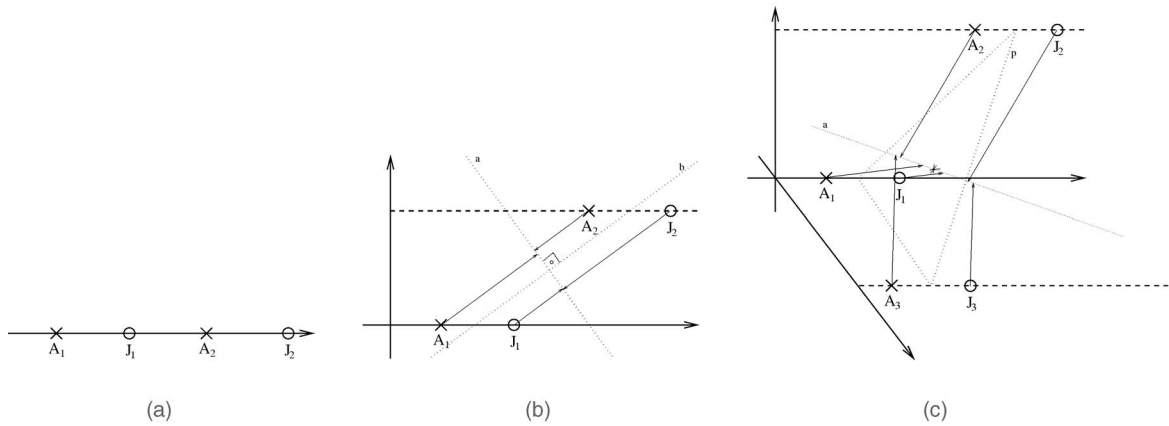


Fig. 3. Illustration of a highly day-dependent feature for two emotions from two different days. (a) The feature values for (A)nger and (J)oy from two different days. (b) Addition of an extra dimension allows for a line b to separate Anger from Joy. The data can be projected down to line a , so the addition of the new dimension did not increase the *final* number of features. (c) In the case of data from three different days, the addition of two extra dimensions allows for a plane p to separate Anger from Joy. The data can again be projected down to line a , not increasing the final number of features.

first factor is slightly more controllable by using disposable electrodes that come from the manufacturer containing a premeasured gel, and always having the subject wash her hands in the same way at the same time before each session, steps we did not impose. However, we made an effort to place the sensors as similarly from day to day as manually possible and to manually apply the same amount of gel each day. Nonetheless, many of these sources of variation are natural and cannot be controlled in realistic long-term measuring applications. Algorithms and features that can compensate for day-to-day variations are needed.

4.3 Day Matrix for Handling Day-Dependence

The 24 statistical features extracted from the signals are dependent on the day the experiment was held. We now augment the 20×24 matrix of the 20 days' 24 features with a 20×19 *day matrix*, which appends a vector of length 19 to each vector of length 24. The vector is the same for all emotions recorded the same day and differs among days. (The vectors were constructed by generating 20 equidistant points in a 19-dim space.) Let us briefly explain the principle with an illustration in Fig. 3.

Consider when the data come from two different days and only one feature is extracted. (This is the simplest way to visualize, but it trivially extends to more features). Although the feature values of one class are always related to the values of the other classes in the same way, e.g., the mean electromyogram signal for anger may always be higher than the mean electromyogram for joy, the actual

values may be highly day-dependent (Fig. 3a). To alleviate this problem, an extra dimension can be added before the features are input into the Fisher Algorithm (Fig. 3b). If the data came from three different days, two extra dimensions are added rather than one (Fig. 3c), etc. In the general case, $D - 1$ extra dimensions are needed for data coming from D different days; hence, we use 19 extra dimensions. The above can be also seen as using the minimum number of dimensions so that each of D points can be at equal distance from all others. Therefore, the $D - 1$ dimensional vector contains the coordinates of one such point for each day.

The effect of the day matrix on classification can be seen in Table 4, where we run Fisher and SFFS-FP with and without the day matrix, and compare it to SFFS, running with the same 24 features as above. Note that it is meaningless to apply SFFS to the day matrix, so that comparison is omitted. The use of the day matrix improves the classification by 3.1, 4.3, 6.9, and 9.4 percent; although only the highest two of these improvements are significant at > 90 percent and > 95 percent (the other confidences are 71-78 percent).

Table 4 also reveals that all the methods perform better on Data Set II than on Data Set I. The improvements range from 5 to 13 percentage points, with confidences 81 percent, 94 percent, 97 percent, 98 percent, and 99 percent.

4.4 Baseline Matrix for Handling Day-Dependence

We propose and evaluate another approach: use of a *baseline matrix* where the Neutral (no emotion) features of each day

TABLE 4
Classification Accuracy for All Eight Emotions for Data Set I and Data Set II Improves with the Use of the Day Matrix

Data Set with 24 Features	Without Day Matrix			With Day Matrix	
	SFFS (%)	Fisher (%)	SFFS-FP (%)	Fisher (%)	SFFS-FP (%)
DataSet I	40.6	40.0	46.3	49.4	50.6
DataSet II	49.4	51.3	56.9	54.4	63.8

These results were obtained with 24 statistical features, $\mu_X, \sigma_X, \delta_X, \tilde{\delta}_X, \gamma_X, \tilde{\gamma}_X, X \in (\mathcal{E}, \mathcal{B}, \mathcal{G}, \mathcal{R})$ for both data sets. The day matrix adds 19 dimensions to each feature that is input into the Fisher Algorithm.

TABLE 5
Data Set I, 7-Emotion (All but the Neutral State)
Classification Results, Comparing Three Methods
for Incorporating the Day Information

Feature Space (Dimensions)	SFFS (%)	Fisher (%)	SFFS-FP (%)
Original (24)	42.9	39.3	45.0
Orig.+ Day (43)	N/A	39.3	45.7
Orig.+ Base. (48)	49.3	40.7	54.3
Orig.+ Base.+ Day (67)	N/A	35.0	49.3

are used as a baseline for (subtracted from) the respective features of the remaining seven emotions of the same day. This gives an additional 20×24 matrix for each of the seven nonneutral emotions. The resulting classification results on seven emotions, run on Data Set I, are in Table 5, together with an additional trial of the day matrix for this case. All the results are significantly higher than random guessing (14.3 percent) with confidence > 99.99 percent. Comparing among the various ways of handling day-dependence, we find the most significant improvement to be SFFS-FP using the baseline matrix, which at 54.3 percent is an improvement over 45 percent (confidence 94 percent). The performance improvements of SFFS-FP over that of Fisher in the last two rows of the table are also statistically significant (confidence 99 percent). Combining all the features (original 24 + baseline + day) appears to result in a decrease in performance, suggesting the limitations of Fisher with too many bad dimensions, where it cannot select out features, but only transform them. However, the decreases (from 54.3 to 49.3 and 40.7 to 35.0) have 80 percent and 84 percent confidence, so are not considered significant.

4.5 Better Features for Handling Day-Dependence

The comparisons above were all conducted with the original six statistical features (1), (2), (3), (4), (5), and (6); we would now like to see how the features $f_1 - f_{10}$ influence classification. We compare four spaces of features from which the algorithms can select and transform subsets: the original six, $\mu_X, \sigma_X, \delta_X, \tilde{\delta}_X, \gamma_X, \tilde{\gamma}_X$, for $X \in (\mathcal{E}, \mathcal{B}, \mathcal{G}, \mathcal{R})$ for a total of 24 features; these same 24 features plus the same six statistics for $X = \mathcal{H}$ for a total of 30 features; features $f_1 - f_{10}$ plus $\mu_{\mathcal{E}}$ which were shown to be useful in our earlier investigation [50], and the combination of all 40 of these.

The results are in Table 6. First, we consider the case where all 40 features are available (the last row). Comparing

each entry in the last row of this table with the same column entry in the first row, we find that all the results are significantly improved (confidence > 99.7 percent).

Next, we look at the results with the day matrix. It appears to improve upon the performance of the statistical features in the first two rows; however, the improvements are only a few percentage points and are not significant (confidences 63-89.5 percent). The additional features of the day matrix can become a liability in the presence of the $f_1 - f_{10}$ features, which compensate much better for the day-to-day variations (the decrease from 70 percent to 61.3 percent is significant at 95 percent confidence.) Without the day matrix, the $f_1 - f_{10}$ features alone perform significantly better than the statistical features (confidence > 93 percent in all six comparisons.)

The best performance occurs when the methods have all forty features at their disposal and the individual best of these occurs with SFFS-FP. The overall best rate of 81.25 percent is significantly higher (> 95 percent confidence) than 16 of the cases shown in Table 6. The three exceptions are in the last row—compared to rates of 78.8 percent and to 77.5 percent the confidences are only 71 percent and 79 percent that 81.25 percent is a genuine improvement.

Table 7 affords a closer look at breaking points of the best-performing case: the SFFS-FP algorithm operating on the full set of 40 features. Summing columns, one sees that the greatest number of false classifications lie in the categories of (P)latoic love and (J)oy. It is possible that this reflects an underlying predilection on the part of this subject toward being in these two categories, even when asked to experience a different emotion; however, one should be careful not to read more into these numbers than is justified.

Because of the longstanding debate regarding whether physiological patterns reveal valence differences (pleasing-displeasing aspects) of emotion, or only arousal differences, we consider here an alternate view of the results in Table 7. Table 8 rearranges the rows and columns into groups based on similar arousal rating or similar valence rating. From this, we see that both arousal and valence are discriminated at rates significantly higher than random (confidence > 99.99 percent), and that the rate of discrimination based on valence (87 percent) versus the rate obtained based on arousal (84 percent) do not differ in a statistically significant way.

TABLE 6
Comparative Classification Rates for Eight Emotions for Data Set II

Number of Features: Which formed initial space	Without Day Matrix			With Day Matrix	
	SFFS (%)	Fisher (%)	SFFS-FP (%)	Fisher (%)	SFFS-FP (%)
24: $\mu_X, \sigma_X, \delta_X, \tilde{\delta}_X, \gamma_X, \tilde{\gamma}_X, X \in (\mathcal{E}, \mathcal{B}, \mathcal{G}, \mathcal{R})$	49.4	51.3	56.9	54.4	63.8
30: $\mu_X, \sigma_X, \delta_X, \tilde{\delta}_X, \gamma_X, \tilde{\gamma}_X, X \in (\mathcal{E}, \mathcal{B}, \mathcal{G}, \mathcal{R}, \mathcal{H})$	52.5	56.9	60.0	58.8	63.8
11: $f_1, f_2, \dots, f_{10}, \mu_{\mathcal{E}}$	60.6	70.0	70.6	61.3	63.1
40: all of the above	65.0	77.5	81.25	77.5	78.8

The day matrix adds 19 features to the data input to the Fisher algorithm, offering no improvement when $f_1 - f_{10}$ are available.

TABLE 7

Confusion Matrix for the Method that Gave the Best Performance Classifying Eight Emotions with Data Set II (81.25 Percent)

	N	A	H	G	P	L	J	R	Total
(N)eutral	17	0	0	0	3	0	0	0	20
(A)nger	0	17	0	0	2	1	0	0	20
(H)atred	0	0	14	1	0	0	3	2	20
(G)rief	0	0	1	15	0	0	4	0	20
(P)latoic Love	0	0	0	0	17	2	1	0	20
Romantic (L)ove	1	1	0	0	3	14	1	0	20
(J)oy	0	0	1	2	0	0	17	0	20
(R)everence	0	0	0	1	0	0	0	19	20
Total	18	18	16	19	25	17	26	21	160

An entry's row label is the true class, the column label is what it was classified as.

4.6 Finding Robust Features

A feature-based approach presents not only a virtually unlimited space of possible features that researchers can propose, but an intractable search of all possible subsets of these to find the best features. When people hand-select features, they may do so with an intuitive feel for which are most important; however, hand-selection is rarely as thorough as machine selection and tends to overlook non-intuitive combinations that may outperform intuitive ones. Machine selection is therefore preferable, even when it is suboptimal (nonexhaustive). Here, we analyze which features the machine searches selected repeatedly. The results of automatic feature selection from twelve experiments involving SFFS (either SFFS as in Jain and Zongker's code, or SFFS-FP, or SFFS-FP with the Day Matrix) are summarized in Table 9.

From Table 9, we see that features such as the means of the heart rate, skin conductivity, and respiration were never selected by any of the classifiers running SFFS, so that they need not have even been computed for these classifiers. At the same time, features such as the mean absolute normalized first difference of the heart rate (δ_H), the first difference of the smoothed skin conductivity (f_4), and the three higher frequency bands of the respiration signal ($f_8 - f_{10}$), were always found to contribute to the best results. Features that were repeatedly selected by classifiers yielding good classification rates can be considered more robust than

those that were rarely selected, but only within the context of these experiments.

There were surprises in that some features physiologists have suggested to be important such as f_3 were not found to contribute to good discrimination. Although the results here will be of interest to psychophysicologists, they must clearly be interpreted in the context of the recognition experiment here and not as some definitive statement about a physical correlate of emotion.

5 CONCLUSIONS

This paper has suggested that machine intelligence should include skills of emotional intelligence, based on recent scientific findings about the role of emotional abilities in human intelligence, and on the way human-machine interaction largely imitates human-human interaction. This is a shift in thinking from machine intelligence as one of primarily mathematical, verbal, and perceptual abilities. Emotion is believed to interact with all of these aspects of intelligence in the human brain. Emotions, largely overlooked in early efforts to develop machine intelligence, are increasingly regarded as an area for important research.

One of the key skills of emotional intelligence for adaptive learning systems is the ability to recognize the emotional communication of others. Even dogs can recognize their owner's affective expressions of pleasure or displeasure—an

TABLE 8

Rearrangements of Table 7, Showing Confusion Matrices for Emotions Having Similar Valence Ratings (139/160 = 87 Percent) and Similar Arousal Ratings (135/160 = 84 Percent)

Valence Rating	Very Neg.	Neg.	Neutral	Pos.	Total
Very Neg. (A)	17	0	0	3	20
Neg. (H,G)	0	31	2	7	40
Neutral (N,R)	0	1	36	3	40
Pos. (P,L,J)	1	3	1	55	60

Arousal Rating	Very High	Med. High	High	Low	Very Low	Total
Very High (A,L)	33	1	0	6	0	40
Med. High (J)	0	17	2	1	0	20
High (G)	0	4	15	1	0	20
Low (N,H,P)	2	4	1	51	2	60
Very Low (R)	0	0	1	0	19	20

TABLE 9

Summary (Along Rows) of when Each Feature Was Chosen in the SFFS-Based Methods of Table 6: The Standard SFFS (S), SFFS Followed by Fisher Projection (SF), and SFFS Followed by Fisher Projection Using the Day Matrix (SFD)

	24 Statistical Features				30 Statistical Features				f_1, f_2, \dots, f_{10} and μ_ε				All 40 Features				Grand	
	S	SF	SFD	Sum	S	SF	SFD	Sum	S	SF	SFD	Sum	S	SF	SFD	Sum	Total	
μ_ε	0	1	1	2/3	0	1	1	2/3	1	1	1	3/3	0	0	1	1/3	8/12	0.67
σ_ε	0	1	1	2/3	0	1	1	2/3					0	1	1	2/3	6/9	0.67
δ_ε	1	1	1	3/3	1	1	1	3/3					1	0	1	2/3	8/9	0.89
$\tilde{\delta}_\varepsilon$	1	1	1	3/3	0	1	1	2/3					0	1	1	2/3	7/9	0.78
γ_ε	1	1	1	3/3	1	1	1	3/3					0	0	1	1/3	7/9	0.78
$\tilde{\gamma}_\varepsilon$	1	1	1	3/3	0	1	1	2/3					0	1	1	2/3	7/9	0.78
μ_B	1	1	1	3/3	0	1	1	2/3					0	1	1	2/3	7/9	0.78
σ_B	0	0	1	1/3	0	0	0	0/3					0	0	0	0/3	1/9	0.11
δ_B	0	0	0	0/3	0	1	1	2/3					0	0	0	0/3	2/9	0.22
$\tilde{\delta}_B$	1	0	1	2/3	0	1	1	2/3					0	1	1	2/3	6/9	0.67
γ_B	0	0	0	0/3	0	0	0	0/3					0	0	0	0/3	0/9	0.00
$\tilde{\gamma}_B$	0	1	1	2/3	0	1	1	2/3					0	0	1	1/3	5/9	0.56
μ_H					0	0	0	0/3					0	0	0	0/3	0/6	0.00
σ_H					0	0	0	0/3					0	0	1	1/3	1/6	0.17
δ_H					0	1	1	2/3					0	1	1	2/3	4/6	0.67
$\tilde{\delta}_H$					1	1	1	3/3					1	1	1	3/3	6/6	1.00
γ_H					0	0	0	0/3					0	1	1	2/3	2/6	0.33
$\tilde{\gamma}_H$					1	1	1	3/3					0	1	1	2/3	5/6	0.83
f_1									1	0	0	1/3	0	0	1	1/3	2/6	0.33
f_2									1	1	1	3/3	0	1	1	2/3	5/6	0.83
μ_S	0	0	0	0/3	0	0	0	0/3					0	0	0	0/3	0/9	0.00
σ_S	0	0	0	0/3	0	0	1	1/3					0	0	1	1/3	2/9	0.22
δ_S	1	1	1	3/3	0	1	1	2/3					1	1	1	3/3	8/9	0.89
$\tilde{\delta}_S$	1	1	1	3/3	0	1	1	2/3					1	1	1	3/3	8/9	0.89
γ_S	1	1	1	3/3	0	1	1	2/3					0	1	1	2/3	7/9	0.78
$\tilde{\gamma}_S$	1	1	1	3/3	0	1	1	2/3					0	1	1	2/3	7/9	0.78
f_3									1	0	0	1/3	0	0	0	0/3	1/6	0.17
f_4									1	1	1	3/3	1	1	1	3/3	6/6	1.00
μ_R	0	0	0	0/3	0	0	0	0/3					0	0	0	0/3	0/9	0.00
σ_R	0	0	1	1/3	0	0	0	0/3					0	0	0	0/3	1/9	0.11
δ_R	1	1	1	3/3	0	1	1	2/3					0	1	1	2/3	7/9	0.78
$\tilde{\delta}_R$	1	1	1	3/3	1	1	1	3/3					0	1	1	2/3	8/9	0.89
γ_R	1	1	1	3/3	0	1	1	2/3					0	1	1	2/3	7/9	0.78
$\tilde{\gamma}_R$	1	1	1	3/3	0	1	1	2/3					0	1	1	2/3	7/9	0.78
f_5									1	0	0	1/3	0	1	1	2/3	3/6	0.50
f_6									1	0	0	1/3	0	1	1	2/3	3/6	0.50
f_7									1	1	1	3/3	0	1	1	2/3	5/6	0.83
f_8									1	1	1	3/3	1	1	1	3/3	6/6	1.00
f_9									1	1	1	3/3	1	1	1	3/3	6/6	1.00
f_{10}									1	1	1	3/3	1	1	1	3/3	6/6	1.00

The features listed in the left-most column are grouped by signal for easier physiological interpretation: electromyogram (ε), blood-volume pressure (B), heart rate (H), skin conductivity (S), and respiration (R). The totals at the right, when high, suggest that the feature may be a robust one regardless of the classification method.

important piece of feedback. One of the difficulties researchers face in this area is the sheer difficulty of getting data corresponding to real emotional states; we have found that efforts in this area are more demanding than traditional efforts to get pattern recognition data. We presented five factors to aid researchers trying to gather good affect data.

The data gathered in this paper contrasts with that gathered for most other efforts at affect pattern recognition not so much in the use of physiology versus video or audio, but in its focus on having the subject try to generate a *feeling*

versus having the subject try to generate an outward *expression*. One weakness of both data-gathering methods is that the subject elicited the emotion, versus a situation or stimulus outside the subject eliciting it. Our group at MIT has recently designed and built environments that focus on emotions not generated deliberately by the subject, e.g., collecting affective data from users driving automobiles in city and highway conditions [38] and from users placed in frustrating computer situations [51]; both of these areas aim at data generation in an *event-elicited*, close to *real-world*,

feeling, open-recording, other-purpose experiment, with effort to make the open recording so comfortable that it is effectively ignored. Nonetheless, much work remains to be done in gathering and analyzing affect data; attempts to create situations that induce emotion in subjects remain subject to uncertainty.

One facet of emotion recognition is developed here for the first time: classification of emotional state from physiological data gathered from one subject over many weeks of data collection. The corpus of person-dependent affect data is larger than any previously reported, and the methodology we developed for its analysis is subject-independent. In future work, the patterns of many individuals may be clustered into groups, according to similarity of the physiological patterns, and these results leveraged to potentially provide person-independent recognition.

Prior to our efforts, researchers have not reported the way in which physiological features for different emotions from the same day tend to be more similar than features for the same emotions on different days. This side-finding of our work may help explain why so many conflicting results exist on prior attempts to identify emotion-dependent physiological correlates. The day-dependent variations that we found and the ways we developed of handling them may potentially be useful in handling across-subject variations. We proposed methods of normalizing features and baselining, showing that the normalizing features gave the best results, but have the drawback of not being easily implemented in an online way because of the requirement of having session-long summarizing information.

We found that the Fisher Projection applied to a subset of features preselected by SFFS always outperformed Fisher Projection applied to the full set of features when assessing percentage of errors made by the classifiers. However, only a few of the improvements of SFFS-FP over FP were significant at > 95 percent; the rest had confidences ranging from 58-87 percent. From the significant improvements, we surmise that the Fisher Projection's ability to transform the feature space may work better when poorer-performing features are omitted up front. The combination is synergistic: fracturization followed by feature transformation, and may well apply in other domains.

Forty features were proposed and systematically evaluated by multiple algorithms. The best and worst features have been identified for this subject. These are of interest in the ongoing search for good features for affect recognition and may aid in trying to understand the differential effects of emotions on physiology.

Although the precise rates found here can only be claimed to apply to one subject, the methodology developed in this paper can be used for any subject. The method is general for finding the features that work best for a given subject and for assessing classification accuracies for that subject.

The results of 81 percent recognition accuracy on eight categories of emotion are the only results we know of for such a classification and are better than machine recognition rates of a similar number of categories of affect from speech (around 60-70 percent) and almost as good as automated recognition of facial expressions (around 80-98 percent).

Because there were doubts in the literature that physiological information shows any differentiation other than arousal level, these results are exciting. However, these results do not imply that a computer can detect and recognize your emotions with 81 percent accuracy because of the limited set of emotions examined, because of the optimistic estimates of error given by the leave-one-out method because of the fact that this experiment, like others in the literature, has only looked at forced choice among presegmented data because of the use of only one subject's long term data and because of the nature of true emotion, which consists of more than externally measurable signals. We expect that joint pattern analysis of signals from face, voice, body, and the surrounding situation is likely to give the most interesting emotion recognition results, especially since people read all these signals jointly.

The greater than six-times-chance recognition accuracy achieved in this pattern recognition research is nonetheless a significant finding; it informs long-debated questions by emotion theorists as to whether there is any physiological differentiation among emotions. Additionally, the finding of significant classification rates when the emotion labels are grouped either by arousal or by valence lends evidence against the belief that physiological signals only differentiate with respect to arousal; both valence and arousal were differentiated essentially equally by the method we developed.

A new question might be stated, "What is the precise bodily nature of the differentiation that exists among different emotions and upon what other factors does this differentiation depend?" We expect that physiological patterning, in combination with facial, vocal, and other behavioral cues, will lead to significant improvements in machine recognition of user emotion over the coming decade, and that this recognition will be critical for giving machines the intelligence to adapt their behavior to interact more smoothly and respectfully with people. However, this "recognition" will be on measurable signals, which we do not expect to include one's innermost thoughts. Much work remains before emotion interpretation can occur at the level of human abilities.

ACKNOWLEDGMENTS

The authors would like to thank Rob and Dan Gruhl for help with hardware for data collection, the MIT Shakespeare Ensemble for acting support, Anil Jain and Doug Zongker for providing us their SFFS code, Tom Minka for help with its use, and the anonymous reviewers for helping improve the paper.

REFERENCES

- [1] A.R. Damasio, *Descartes' Error: Emotion, Reason, and the Human Brain*. New York: Gosset/Putnam Press, 1994.
- [2] J. LeDoux, *The Emotional Brain*. New York: Simon & Schuster, 1996.
- [3] P. Salovey and J.D. Mayer, "Emotional Intelligence," *Imagination, Cognition and Personality*, vol. 9, no. 3, pp. 185-211, 1990.
- [4] D. Goleman, *Emotional Intelligence*. New York: Bantam Books, 1995.
- [5] B. Reeves and C. Nass, *The Media Equation*. Cambridge Univ. Press, Center for the Study of Language and Information, 1996.

- [6] M. Sigman and L. Capps, *Children with Autism: A Developmental Perspective*. Cambridge, Mass.: Harvard Univ. Press, 1997.
- [7] R.W. Picard, *Affective Computing*. Cambridge, Mass.: The MIT Press, 1997.
- [8] J. Healey, R.W. Picard, and F. Dabek, "A New Affect-Perceiving Interface and Its Application to Personalized Music Selection," *Proc. Workshop Perceptual User Interfaces*, Nov. 1998.
- [9] R.W. Picard and J. Healey, "Affective Wearables," *Personal Technologies*, vol. 1, no. 4, pp. 231-240, 1997.
- [10] K.R. Scherer, "Speech and Emotional States," *Speech Evaluation in Psychiatry*, J.K. Darby, ed., chapter 10, pp. 189-220, Grune and Stratton, Inc., 1981.
- [11] R. Banse and K.R. Scherer, "Acoustic Profiles in Vocal Emotion Expression," *J. Personality and Social Psychology*, vol. 70, no. 3, pp. 614-636, 1996.
- [12] T. Polzin, "Detecting Verbal and Non-Verbal Cues in the Communication of Emotions," PhD thesis, School of Computer Science, June 2000.
- [13] J. Hansen, Comm. during ICASSP '99 Panel on Speech Under Stress. *Proc. Int'l Conf. Acoustics, Speech, and Signal Processing '99*, 1999.
- [14] J.N. Bassili, "Emotion Recognition: The Role of Facial Movement and the Relative Importance of Upper and Lower Areas of the Face," *J. Personality and Social Psychology*, vol. 37, pp. 2049-2058, 1979.
- [15] Y. Yacoob and L.S. Davis, "Recognizing Human Facial Expressions from Log Image Sequences Using Optical Flow," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 18, no. 6, pp. 636-642, June 1996.
- [16] I. Essa and A. Gardner, "Prosody Analysis for Speaker Affect Determination," *Proc. Workshop Perceptual User Interfaces '97*, pp. 45-46, Oct. 1997.
- [17] I.A. Essa, "Analysis, Interpretation and Synthesis of Facial Expressions," PhD thesis, Mass. Inst. of Technology, Media Lab, Cambridge, Mass., Feb. 1995.
- [18] J.F. Cohn, A.J. Zlochower, J. Lien, and T. Kanade, "Automated Face Analysis by Feature Point Tracking has High Concurrent Validity with Manual FACS Coding," *Psychophysiology*, vol. 36, pp. 35-43, 1999.
- [19] M. Bartlett, J.C. Hager, P. Ekman, and T.J. Sejnowski, "Measuring Facial Expressions by Computer Image Analysis," *Psychophysiology*, vol. 36, pp. 253-263, 1999.
- [20] G. Donato, M.S. Bartlett, J.C. Hager, P. Ekman, and T.J. Sejnowski, "Classifying Facial Actions," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 21, no. 10, pp. 974-989, Oct. 1999.
- [21] L.C. DeSilva, T. Miyasato, and R. Nakatsu, "Facial Emotion Recognition Using Multi-Modal Information," *Proc. IEEE Int'l Conf. Information, Comm., and Signal Processing*, pp. 397-401, Sept. 1997.
- [22] T.S. Huang, L.S. Chen, and H. Tao, "Bimodal Emotion Recognition by Man and Machine," *Proc. ATR Workshop Virtual Communication Environments*, Apr. 1998.
- [23] L.S. Chen, T.S. Huang, T. Miyasato, and R. Nakatsu, "Multimodal Human Emotion/Expression Recognition," *Proc. Third Int'l Conf. Automatic Face and Gesture Recognition*, pp. 366-371, Apr. 1998.
- [24] J.T. Cacioppo and L.G. Tassinari, "Inferring Psychological Significance from Physiological Signals," *Am. Psychologist*, vol. 45, pp. 16-28, Jan. 1990.
- [25] W. James, *William James: Writings 1878-1899*, pp. 350-365, The Library of Am., 1992. Originally published in 1890.
- [26] W.B. Cannon, "The James-Lange Theory of Emotions: A Critical Examination and an Alternative Theory," *Am. J. Psychology*, vol. 39, pp. 106-124, 1927.
- [27] S. Schachter, "The Interaction of Cognitive and Physiological Determinants of Emotional State," *Advances in Experimental Psychology*, L. Berkowitz, ed., vol. 1, pp. 49-80, 1964.
- [28] P. Ekman, R.W. Levenson, and W.V. Friesen, "Autonomic Nervous System Activity Distinguishes Among Emotions," *Science*, vol. 221, pp. 1208-1210, Sept. 1983.
- [29] W.M. Winton, L. Putnam, and R. Krauss, "Facial and Autonomic Manifestations of the Dimensional Structure of Emotion," *J. Experimental Social Psychology*, vol. 20, pp. 195-216, 1984.
- [30] A.J. Fridlund and C.E. Izard, "Electromyographic Studies of Facial Expressions of Emotions and Patterns of Emotions," *Social Psychophysiology: A Sourcebook*, J.T. Cacioppo and R.E. Petty, eds., pp. 243-286, 1983.
- [31] J.T. Cacioppo, G.G. Berntson, J.T. Larsen, K.M. Poehlmann, and T.A. Ito, "The Psychophysiology of Emotion," *Handbook of Emotions*, M. Lewis and J.M. Haviland-Jones, eds., pp. 173-191, 2000.
- [32] D. Keltner and P. Ekman, "The Psychophysiology of Emotion," *Handbook of Emotions*, M. Lewis and J.M. Haviland-Jones, eds., pp. 236-249, 2000.
- [33] D.M. Clynes, *Sentics: The Touch of the Emotions*. Anchor Press/Doubleday, 1977.
- [34] C.E. Izard, "Four Systems for Emotion Activation: Cognitive and Noncognitive Processes," *Psychological Rev.*, vol. 100, no. 1, pp. 68-90, 1993.
- [35] H. Hama and K. Tsuda, "Finger-Pressure Waveforms Measured on Clynes' Sentograph Distinguished Among Emotions," *Perceptual and Motor Skills*, vol. 70, pp. 371-376, 1990.
- [36] H. Schlosberg, "Three Dimensions of Emotion," *Psychological Rev.*, vol. 61, pp. 81-88, Mar. 1954.
- [37] P.J. Lang, "The Emotion Probe: Studies of Motivation and Attention," *Am. Psychologist*, vol. 50, no. 5, pp. 372-385, 1995.
- [38] J.A. Healey, "Wearable and Automotive Systems for Affect Recognition from Physiology," Technical Report 526, PhD thesis, Mass. Inst. Technology, Cambridge, Mass., May 2000.
- [39] E. Vyzas and R.W. Picard, "Affective Pattern Classification," *Proc. AAAI 1998 Fall Symp., Emotional and Intelligent: The Tangled Knot of Cognition*, Oct. 1998.
- [40] E. Vyzas and R.W. Picard, "Offline and Online Recognition of Emotion Expression from Physiological Data," *Proc. Workshop Emotion-Based Agent Architectures, Third Int'l Conf. Autonomous Agents*, pp. 135-142, May 1999.
- [41] A.V. Oppenheim and R.W. Schaefer, *Discrete-Time Signal Processing*. Englewood Cliffs, N.J.: Prentice-Hall, 1989.
- [42] D.T. Lyyken, R. Rose, B. Luther, and M. Maley, "Correcting Psychophysiological Measures for Individual Differences in Range," *Psychophysiological Bulletin*, vol. 66, pp. 481-484, 1966.
- [43] D.T. Lyyken and P.H. Venables, "Direct Measurement of Skin Conductance: A Proposal for Standardization," *Psychophysiology*, vol. 8, no. 5, pp. 656-672, 1971.
- [44] M.E. Dawson, A.M. Schell, and D.L. Filion, "The Electrodermal System," *Principles of Psychophysiology: Physical, Social, and Inferential Elements*, J.T. Cacioppo and L.G. Tassinari, eds., pp. 295-324, 1990.
- [45] P. Pudil, J. Novovicova, and J. Kittler, "Floating Search Methods in Feature Selection," *Pattern Recognition Letters*, vol. 15, pp. 1119-1125, Nov. 1994.
- [46] A. Jain and D. Zongker, "Feature Selection: Evaluation, Application, and Small Sample Performance," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 19, no. 2, pp. 153-163, Feb. 1997.
- [47] P. Pudil, J. Novovicova, and J. Kittler, "Simultaneous Learning of Decision Rules and Important Attributes for Classification Problems in Image Analysis," *Image and Vision Computing*, vol. 12, pp. 193-198, Apr. 1994.
- [48] R.O. Duda and P.E. Hart, *Pattern Classification and Scene Analysis*. Wiley-Interscience, 1973.
- [49] T.M. Cover and P.E. Hart, "Nearest Neighbor Pattern Classification," *IEEE Trans. Information Theory*, vol. 13, pp. 21-27, Jan. 1967.
- [50] J. Healey and R.W. Picard, "Digital Processing of Affective Signals," *Proc. IEEE Int'l Conf. Acoustics, Speech, and Signal Processing*, 1998.
- [51] J. Scheirer, J. Klein, R. Fernandez, and R.W. Picard, "Frustrating the User on Purpose: A Step Toward Building an Affective Computer," *Interaction with Computers*, 2001, to appear.



Rosalind W. Picard earned the bachelors degree in electrical engineering with highest honors from the Georgia Institute of Technology in 1984. She was named a National Science Foundation Graduate Fellow and worked as a member of the technical staff at AT&T Bell Laboratories from 1984-1987, designing VLSI chips for digital signal processing and developing new methods of adaptive image compression. She earned the masters and

doctorate degrees, both in electrical engineering and computer science, from the Massachusetts Institute of Technology (MIT) in 1986 and 1991, respectively. In 1991, she joined the MIT Media Laboratory as an assistant professor. She was promoted to associate professor in 1995 and awarded tenure at MIT in 1998. She is author or co-author of more than 80 peer reviewed scientific articles in pattern recognition, multi-dimensional signal modeling, computer vision, and human-computer interaction. She is a co-recipient with Tom Minka of a best paper prize (1998) from the Pattern Recognition Society for work on machine learning with multiple models. Dr. Picard guest edited the *IEEE Transactions on Pattern Analysis and Machine Intelligence* special issue on Digital Libraries: Representation and Retrieval, and edited the proceedings of the First IEEE International Workshop on Content-Based Access of Image and Video Libraries, for which she served as chair. Her award-winning book, *Affective Computing*, (MIT Press, 1997) lays the groundwork for giving machines skills of emotional intelligence. She is a senior member of the IEEE and a member of the IEEE Computer Society.



Elias Vyzas received the BEngr degree in mechanical engineering from Imperial College, London, England in June 1994, and two MS degrees, one in mechanical engineering and one in electrical engineering and computer science, from MIT in January 1997. He received the masters degree in mechanical engineering from MIT in June 1999 and is currently completing service in the Greek Army. His research interests include affect and stress recognition from physiology, and assessment of mental workload and performance.



gence and knowledge management. She is currently a research staff member in the Human Language Technologies Group at IBM's T.J. Watson Research Center.

Jennifer Healey received the BS (1993), the MS (1995) and the PhD (2000) degrees from MIT in electrical engineering and computer science. Her Master's thesis work in the field of optics and opto-electronics was completed at the Charles Stark Draper Laboratory. In 1995, she joined the MIT Media Laboratory and helped launch the Affective Computing research project. She held a postdoctoral position at IBM's Zurich Research Laboratory in artificial intelli-

► **For more information on this or any other computing topic, please visit our Digital Library at <http://computer.org/publications/dlib>.**