

Robust Recognition of Emotion from Speech

Mohammed E. Hoque¹, Mohammed Yeasin¹, and Max M. Louwerse²

¹ Department of Electrical and Computer Engineering / Institute for Intelligent Systems

² Department of Psychology / Institute for Intelligent Systems

The University of Memphis

Memphis, TN 38152 USA

{mhoque, myeasin, mlouwerse}@memphis.edu

Abstract. This paper presents robust recognition of a subset of emotions by animated agents from salient spoken words. To develop and evaluate the model for each emotion from the chosen subset, both the prosodic and acoustic features were used to extract the intonational patterns and correlates of emotion from speech samples. The computed features were projected using a combination of linear projection techniques for compact and clustered representation of features. The projected features were used to build models of emotions using a set of classifiers organized in hierarchical fashion. The performances of the models were obtained using number of classifiers from the WEKA machine learning toolbox. Empirical analysis indicated that the lexical information computed from both the prosodic and acoustic features at word level yielded robust classification of emotions.

Keywords: emotion recognition, prosody, speech, machine learning.

1 Introduction

Animated conversational agents allow for natural multimodal human-computer interaction and have shown to be effective in various intelligent systems, including intelligent tutoring systems [1, 2]. Agents used in intelligent tutoring are designed to articulate difficult concepts in a well paced, adaptive and responsive atmosphere based on the learners' affective and cognitive states. Expert educators, both human and artificial, are expected to identify the cognitive states of mind of the learners' and take appropriate pedagogical actions [3]. Because of the realization that monitoring cognitive states in the student through the student's verbal feedback alone is not enough, research that focuses on monitoring of other modalities like speech has become more common [4, 5]. There is no doubt that high accuracy recognition of cognitive states and emotions relies on multiple modalities, rather than one specific modality. For instance, when a speaker is surprised, this emotion can be expressed through language (syntax), facial expressions (eyebrows moving up), through gestures (showing palms of both hands), through eye gaze (making continued eye contact with dialogue partner) as well as through speech (intonational contours). Moreover, in addition to multiple modalities being responsible for the recognition of emotions and

cognitive states, the interaction of modalities is of importance, because one modality can compensate the absence of another. Because of the current state of human-computer interaction, the modalities animated conversational agents can use for their response to a dialogue partner are speech and language.

Despite the fact that we know linguistic modalities (e.g. dialog move, intonation, pause) and paralinguistic modalities (e.g. facial expressions, eye gaze, gestures) interact in communication, the exact nature of this interaction remains unclear [6]. There are two primary reasons why an insight in the interaction of modalities in the communicative process is beneficial.

First, from a psychological point of view it helps us understand how communicative processes take shape in the minds of dialog participants. Under what psychological conditions are different channels aligned? Does a channel add information to the communicative process or does it merely co-occur with other channels? Research in psychology has shed light on the interaction of modalities, for instance comparing eye gaze [7, 8], gestures [9-11] and facial expressions [12] but many questions regarding multiple – i.e., more than pairs of – channels and their alignment remain unanswered.

Second, insight in multimodal communication is beneficial from a computational point of view, for instance in the development of animated conversational agents [13]. The naturalness of the human-computer interaction can be maximized by the use of animated conversational agents, because of the availability of both linguistic (semantics, syntax) and paralinguistic (pragmatic, sociological) features. These animated agents have anthropomorphic, automated, talking heads with facial features and gestures that are coordinated with text-to-speech-engines [14-16]). Examples of these agents are Baldi [17], COSMO [18], STEVE [19], Herman the Bug [18] and AutoTutor [20]. Though the naturalness of these agents is progressively changing, there is room for improvement. Current agents for instance incessantly stare at the dialog partner, use limited facial features rather randomly, or produce bursts of unpaused speech. Both psycholinguistics and computational linguistics would thus benefit from answers to questions regarding the interaction of multimodal channels.

There is a growing interest in robust recognition of emotion from speech by researchers from various interdisciplinary areas. Examples of specific domains are affective interface [3] and call center environments [21]. In recent work by Dellaert *et al.* [22] accuracies in the range of 60% -65% were reported in distinguishing patterns among sadness, anger, happiness, and fear in the general domain of Human-Computer Interaction (HCI). The results were obtained using a cross-validation approach by fusing three classifiers: the maximum likelihood Bayes classification, kernel regression, and the k -nearest neighbor (k -NN) methods using the pitch contour features. For a call center environment Lee *et al.* [23] distinguish between two emotions: positive and negative, using linear discrimination, k -NN classifiers, and support vector machines achieving a maximum accuracy rate of 75%.

Paeschke [24] used a real-time emotion recognizer with neural networks adopted for call center applications and reported 77% classification accuracy in two emotions: agitation and calmness. Several studies showed how “quality features” (based on formant analysis) are used in addition to “prosody features”, (particularly pitch and

energy) to improve the classification of multiple emotions [25], [26]. This technique is known to exploit emotional dimensions other than prosody. Yu *et al.* [27] used SVMs, binary classifiers, to detect one emotion versus the rest. On four distinct emotions such as anger, happiness, sadness, and neutral, they achieved an accuracy of 73%.

Robust recognition of emotion expressed in speech requires a thorough understanding of the lexical aspects of speech [21]. Lee *et al.* hypothesized that a group of positive and negative words were related to different emotions. The occurrences of such predefined words were used to infer the emotional reaction of a caller using a probabilistic framework. Lee *et al.* argued that there a one-to-one correspondence may be assumed between a word and a positive or negative emotion.

Though this may be true for some words that have a semantic bias, more commonly a word does not have such a bias and can convey different emotions by the use of different intonational patterns. For example, the frequently used discourse marker “okay”, is often used to express affirmation (S1 “Ready?” S2 “Okay”), but can also be used to express delight (S1 “So and that’s how the procedure works” S2 “Okay!”), confidence (S1 “You’re ready for the jump?” “Okay”), or confusion (S1 “You just multiply by the divider” S2 “Okay...?”) [28]. The meaning of these different uses of “okay” may be guessed by their context, but to a large extent their

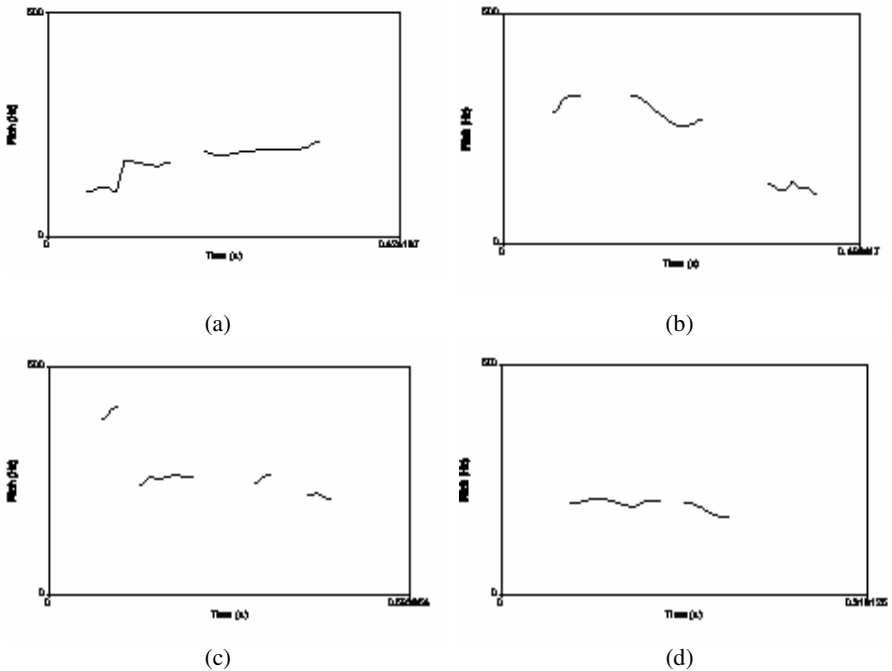


Fig. 1. Pictorial depiction of the word “okay” uttered with different intonations to express different emotions. The pitch contour of various emotions: (a) confusion, (b) flow, (c) delight and (d) neutral are plotted to highlight the differences manifested at lexical level by various emotions.

emotional value only becomes clear in the intonational patterns used to express the word. Figure 1 shows that despite the fact that a word like “okay” is the same, the intonational patterns are very different depending on the emotions. We therefore predict that lexical information extracted from combined prosodic and acoustic features that correspond to intonational pattern of “salient words” will yield robust recognition of emotion from speech, providing a framework for signal level analysis of speech for emotion.

To test this hypothesis, a small database of audio samples representing various emotions was used. Based on the domain knowledge, preprocessing of audio samples is performed to extract the salient words and selected word-level utterances were used to compute features such as fundamental frequency (F0), energy, rhythm, pause and duration. The computed features were projected and then, fused in a feature level framework to build models for various emotions.

2 Proposed Approach

The proposed approach consists of five major components (see Figure 2): (i) collection of suitable data sets for training and testing, (ii) extraction of feature, (iii) projection of feature to lower dimensional space, (iv) learning the models using machine learning techniques and (v) evaluation of models. This paper thereby presents a holistic approach in robust recognition of emotion from speech.

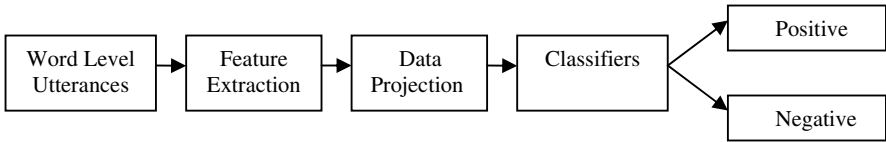


Fig. 2. The high level description of the overall emotion recognition process

First, a suitable database is captured for building and evaluating the models. Second, intonational patterns from spoken “salient words” are extracted with a combination of prosodic and acoustic features. Third, the extracted features are projected onto the lower dimensional space using combined Principle Component Analysis (PCA) [29] and Linear Discriminant Analysis (LDA) for a compact and clustered representation of computed features. Fourth, a set of machine learning techniques from the WEKA [30] toolbox are used to learn the models from the training samples. Finally, testing samples are used to evaluate the performance of the models. We describe the details of various components of robust recognition of emotion from speech below.

2.1 Database and Preparation

Collecting large databases of natural and unbiased emotion is challenging. One needs a representative data set to infer various emotions from speech using machine learning technique to establish the hypothesis and to obtain meaningful results. The

performance of a classifier that can distinguish different emotional patterns ultimately depends on the completeness of the training and testing samples and how similar these samples are to real-world data.

The data captured to perform experiments can be categorized into three methods depending on how they are captured. The first method employs actors to utter various or similar sentences in various feigned emotional patterns. The second method utilizes a system that interacts with a human subject and draws him/her to an emotional point and records the response. The third approach is to extract real-life human utterances, which express various natural emotions.

The main drawback of having actors expressing emotional utterance is that the utterances are acted out independently from one another typically in a laboratory setting. These data may converge very well, but may not be suitable for real-life human-computer interaction settings. On the other hand, setting up an experiment where individuals interact with computers or other individuals is expensive and time consuming for testing out classifiers. In the study reported here, emotional utterances were clipped from movies. Though it is true that emotions are still “acted out”, the discourse context and the absence of a lab setting makes it more natural than the first method. Utterances were taken from three movies: “Fahrenheit 911”, “Bowling for Columbine” and “Before Sunset”. “Fahrenheit 911” and “Bowling for Columbine” are political documentaries containing real interviews with many cases of positive and negative emotions. “Before Sunset” is a chatty romantic movie with delightful, frustrating and confusing expressions with minimal background music. Fifteen utterances were selected from these movies covering four classes of emotions: confusion/uncertain, delight, flow (confident, encouragement), and frustration [3, [4, 5]. Selected utterances were stand-alone expressions in conversations that had an ambiguous meaning, dependent on the context (e.g. “Great”, “Yes”, “Yeah”, “No”, “Ok”, “Good”, “Right”, “Really”, “What”, “God”). Three graduate students listened to the audio clips without specific instructions as to what intonational patterns to listen to and successfully distinguished between the positive and negative emotions 65% of the time. A hierarchical classifier was designed to first distinguish between positive (delight and flow) and negative (confusion and frustration) emotions. The same set of classifiers were applied again on positive and negative emotions

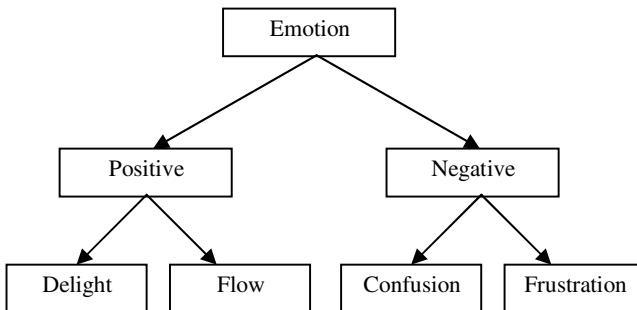


Fig. 3. The design of the hierarchical binary classifiers

separately to differentiate between delight and flow under positive emotions, and confusion and frustration under negative emotions as shown in Figure 3.

2.2 Emotion Models Using Lexical Information

To compute the lexical information from spoken salient words, 22 acoustic and prosodic features related to segmental and suprasegmental information, which are believed to be correlates of emotion, were calculated. Computed features were utterance level statistics related to fundamental frequency (F0) [31-33]. Other features were related to duration, intensity, and formants. In particular, the following features were computed for developing the models.

1. **Pitch:** Minimum, maximum, mean, standard deviation, absolute value, quantile, ratio between voiced and unvoiced frames.
2. **Duration:** ϵ_{time} ϵ_{height}
3. **Intensity:** Minimum, maximum, mean, standard deviation, quantile.
4. **Formant:** First formant, second formant, third formant, fourth formant, fifth formant, second formant / first formant, third formant / first formant
5. **Rhythm:** Speaking rate.

The speech processing software Praat [34] was used to calculate the features in batch mode. ϵ_{time} , ϵ_{height} features, which are part of duration, are prominent measures.

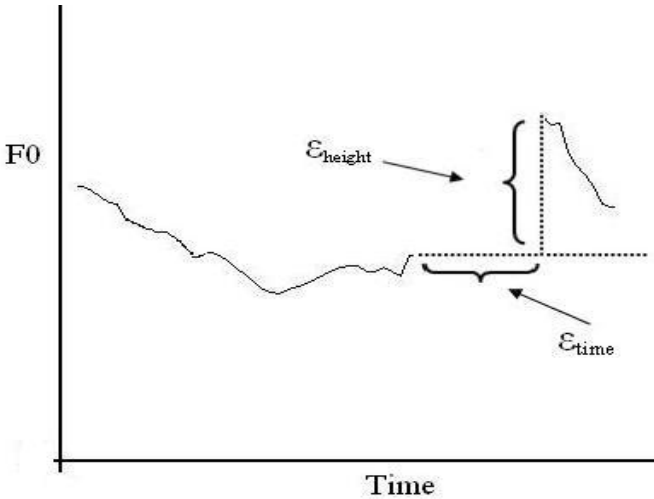


Fig. 4. Measures of F0 for computing parameters (ϵ_{time} , ϵ_{height}) which corresponds to rising and lowering of intonational

ϵ_{height} and ϵ_{time} features are related to phenomenon when fundamental frequency breaks down in word levels. ϵ_{time} refers to the pause time between two disjoint segments of F0 (often referred as Pitch), whereas ϵ_{height} refers to the vertical distance between the segments symbolizing voice breaks as shown in Figure 4. Inclusion of

height and *time* accounts for possible low or high pitch accents. The frequency shift between the segments was selected rather than absolute measures to take into account the discourse [35].

The first model fed the raw 22 features directly into the classifier. The second and the third model applied PCA on the raw features and took the first 15 (F15) and 20 (F20) eigenvectors respectively to de-correlate the base features. In the fourth model, LDA is directly used on the raw features to project them directly onto lower dimension. The fifth model consisted of the combination of PCA (F15) and LDA. A 10-fold cross validation technique was used whereby the training data was randomly split into ten sets, 9 of which were used in training and the 10th for validations. Then iteratively another nine were picked.

Table 1. The list of classifiers used to validate the robustness of the algorithm using weka toolbox

Types of Classifiers				
Rules	Trees	Meta	Functions	Bayes
Part	RandomForrest	AdaBoostM1	Logistic	Naïve Bayes
NNge	J48	Bagging	Multi-layer Perceptron	Naïve Bayes Simple
Ridor	Logistic Model Tree	Classification via Regression	RBF Network	Naïve Bayes Updateable
-	-	LogitBoost	Simple Logistics	-
-	-	Multi Class Classifier	SMO	-
-	-	Ordinal Class Classifier	-	-
-	-	Threshold Selector	-	-

2.3 Results and Discussion

Results showed that the combination of data projection techniques such as PCA and LDA yielded better performance as opposed to using raw features or using LDA or PCA alone (Table 2). An average of 83.33% accuracy was achieved using the combination of PCA and LDA. On the other hand, features like PCA (F15), PCA (F20) and LDA resulted in accuracy rates of respectively 50.79%, 57.1%, 61%, and 52.01% on average. The performance of combining PCA and LDA is higher than PCA or LDA itself mainly because PCA de-correlates the data, whereas LDA projects the data onto lower dimension. Therefore, the combination of PCA and LDA is expected to work better.

When the same models were applied to positive emotions and negative emotions separately even more impressive results emerged (Table 3). The performance of the diverse set of classifiers to recognize negative emotions is better than the performance to recognize positive emotions. One potential explanation for this is that negative

emotions may deviate more from the standard than positive emotions. In other words, positive emotions may in general less likely be recognized as an emotion, because they map onto the default. Negative emotions on the other hand deviate from that default, thereby facilitating recognition, both in humans and computers.

Table 2. Summary of classification results for 21 selected classifiers

Category	Classifiers	Accuracy (%)				
		Features (a)	PCA (b)		LDA (c)	PCA+LDA (d)
			F15 (b1)	F20 (b2)		
Rules	Part	50	66.67	66.67	47.61	83.33
	NNge	33.33	33.33	38.09	38.09	83.33
	Ridor	66.67	83.33	100	47.20	66.67
Trees	Random Forrest	50	50	50	66.67	83.33
	J48	50	66.67	66.67	47.61	83.33
	Logistic Model Tree	33.33	47.61	83.33	66.67	71.67
Meta	AdaBoostM1	61.90	71.42	71.42	42.85	61.90
	Bagging	33.33	66.67	83.33	42.85	66.67
	Classification via Regression	50	66.67	66.67	47.61	83.33
	Logit Boost	50	50	61.90	52.38	83.33
	Multi Class Classifier	50	42.85	52.38	57.14	83.33
	Ordinal Class Classifier	50	66.67	66.67	47.62	83.33
	Threshold Selector	50	66.67	66.67	61.90	100
Functions	Logistic	50	42.85	57.38	57.14	83.33
	Multi-layer Perceptron	50	57.14	52.38	50	83.33
	RBF Network	33.33	66.67	52.38	38.09	83.33
	Simple Logistics	33.33	47.61	83.33	66.67	66.67
	SMO	71.42	57.14	61.90	52.38	71.42
Bayes	Naïve Bayes	66.67	50	33.33	52.38	66.67
	Naïve Bayes Simple	66.67	50	33.33	57.14	66.67
	Naïve Bayes Updateable	66.67	50	33.33	52.38	66.67

Note. (a) raw features are used into classifiers, (b1) using the first 15 (f15) eigenvectors of PCA into the classifiers, (b2) using the first 20 (f20) eigenvectors of PCA into the classifiers. (c) using LDA to project the data into lower dimension and then use them into the classifiers. (d) combination of both PCA and LDA to not only de-correlate the data redundant feature space, but also to project them into lower dimension and then use them into the classifiers.

Table 3. Summary of classification results for 21 classifiers on positive and negative emotions

Category	Classifiers	Accuracy (%)	
		Delight + Flow	Confusion + Frustration
Rules	Part	72.72	100
	NNge	80	100
	Ridor	66.67	100
Trees	RandomForrest	63.63	66.67
	J48	72.72	100
	LMT	72.72	100
Meta	AdaBoostM1	54.44	100
	Bagging	63.64	66.67
	Classification via Regression	72.72	100
	LogitBoost	63.64	100
	Multi Class Classifier	72.72	100
	Ordinal Class Classifier	72.72	100
	Threshold Selector	83.33	100
Functions	Logistic	72.72	100
	Multi-layer Perceptron	66.67	100
	RBF Network	66.67	100
	Simple Logistics	72.72	100
	SMO	72.72	100
Bayes	Naïve Bayes	72.72	100
	Naïve Bayes Simple	72.72	100
	Native Bayes Updateable	72.72	100

Note. Results with the combination of PCA + LDA were only recorded as they comparatively produce better results as shown in Table 2.

It needs to be noted that the results presented in Table 2 and 3 are satisfactory but very similar. The most likely explanation for this result is the limited dataset that does not provide the variability that can be found in larger sets of spoken discourse. Indeed, additional data collection needs to be conducted in order to provide a larger sample set for training and testing.

The total classifiers used can be broken into five categories. Rule based classifiers produce rules for classification from the training data and then apply them on the testing set. Tree based classifiers produce classification trees as their outputs. Function-based classifiers, on the other hand, represent the well-known support vector machine, neural network, linear regression types of classifiers. Meta classifiers combine several classifiers, e.g. Vote or enhance a single classifier, e.g. bagging. Bayes group consists of simple probabilistic classifiers. From Tables 2 & 3, it is evident that, with the exception of Bayes, all the classifiers perform similarly in this particular problem domain. It can be easily explained that due to the limited database, probabilistic based classifiers such as Bayes did not perform equally well compared to the other classifiers. In the second

phase of this study with more challenging map-task data, similar performances across a variety of classifiers would be unlikely. This may provide a better conclusive result about a set of optimum classifiers for this given problem.

3 Conclusions

Robust autonomous recognition of emotion is gaining attention due to the widespread applications into various domains, including those with animated conversational agents. Automated recognizing emotion with high accuracy still remains an elusive goal due to the lack of complete understanding and agreement of emotion in human minds. The study presented in this paper achieved an average of 83.33% success rate of defining positive and negative emotion using a varied set of classifiers confined to learning environment. Lexical and prosodic features were used on word level emotional utterances to improve the performance of the emotion recognition system. Our results indicate that using a proper set of projection techniques on word level lexical and prosodic features yields an accuracy rate of 80 to 100%. It is worth noting that the datasets were tested by three graduate students who were able to classify the emotions into correct bins 65% of the time. This supports our hypothesis that word level prosodic and lexical features provide useful clues about positive and negative emotions. This hypothesis also enables us to have a framework for signal level analysis.

Obviously, there is a risk involved in clipping arbitrary words from a conversation, which may be ineffective at various cases as some words may convey more in context only. Therefore, our goal for the immediate future is to look at meaningful words in a sequence while introducing context in our analysis as well. A research project that investigates multimodal communication (prosody, dialog structure, eye gaze and facial expressions) in Map Task scenarios will thereby generate the needed data [5, 6]. In the second phase of this project the results of the data analysis will allow us to develop an animated conversational agent that uses the right intonational contours in the right contexts, expressing the right emotions.

Visual information modifies the perception of speech [17], while combinations of visual and audio information provide robust performance when modalities are captured in noisy environment [36]. Therefore, in order for an animated conversational agent to be successful in learning environments, it is imperative that the agent should be able to fuse the audio and video data to reach a decision regarding the emotional states of the learners. Therefore, our future efforts will include fusion of video and audio data in a signal level framework to boost the performance of our existing emotion recognition system.

Acknowledgements

This research was partially supported by grant NSF-IIS-0416128 awarded to the third author. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the funding institution.

References

- [1] A. C. Graesser, K. VanLehn, C. Rose, P. Jordan, and D. Harter, "Intelligent tutoring systems with conversational dialogue.," *AI Magazine*, vol. 22, pp. 39-51, 2001.
- [2] M. M. Louwerse, A. C. Graesser, S. Lu, and H. H. Mitchell, "Social cues in animated conversational agents," *Applied Cognitive Psychology*, vol. 19, pp. 1-12, 2005.
- [3] B. Kort, R. Reilly, and R. W. Picard, "An Affective Model of Interplay Between Emotions and Learning: Reengineering Educational Pedagogy-Building a Learning Companion.," presented at In Proceedings of International Conference on Advanced Learning Technologies (ICALT 2001), Madison, Wisconsin, August 2001.
- [4] S. K. D'Mello, S. D. Craig, A. Witherspoon, J. Sullins, B. McDaniel, B. Gholson, and A. C. Graesser, "The relationship between affective states and dialog patterns during interactions with AutoTutor," presented at Proceedings of the World Conference on E-learning in Corporate, Government, Health Care, and Higher Education, Chesapeake, VA, 2005.
- [5] M. Louwerse, P. Jeuniaux, M. Hoque, J. Wu, and G. Lewis, "Multimodal Communication in Computer-Mediated Map Task Scenarios," presented at The 28th Annual Conference of the Cognitive Science Society, Vancouver, Canada, July 2006.
- [6] M. M. Louwerse, E. G. Bard, M. Steedman, X. Hu, and A. C. Graesser, "Tracking multimodal communication in humans and agents," Institute for Intelligent Systems, University of Memphis, Memphis, TN., 2004.
- [7] M. Argyle and M. Cook, *Gaze and Mutual Gaze*, 1976.
- [8] G. Doherty-Sneddon, A. H. Anderson, C. O'Malley, S. Langton, S. Garrod, and V. Bruce, "Face-to-face and video-mediated communication: A comparison of dialogue structure and task performance," *Journal of Experimental Psychology-Applied*, vol. 3(2), pp. 105-125, 1997.
- [9] S. Goldin-Meadow and M. A. Singer, "From children's hands to adults' ears: Gesture's role in teaching and learning.," *Developmental Psychology*, pp. 509-520, 2003.
- [10] M. M. Louwerse and A. Bangerter, "Focusing attention with deictic gestures and linguistic expressions," presented at Proceedings of the Cognitive Science Society, 2005.
- [11] D. McNeill, "Hand and mind: What gestures reveal about thought.," 1992.
- [12] P. Ekman, "About brows: emotional and conversational signals.," *Human Ethology*, pp. 169-248, 1979.
- [13] M. M. Louwerse, A. C. Graesser, S. Lu, and H. H. Mitchell, "Social cues in animated conversational agents," *Applied Cognitive Psychology*, vol. 19, pp. 1-12, 2004.
- [14] J. Cassell and K. Thorisson, "The Power of a Nod and a Glance: Envelope vs. Emotional Feedback in Animated Conversational Agents.," *Journal of Applied Artificial Intelligence*, vol. 13 (3), pp. 519-538, 1999.
- [15] M. M. Cohen and D. W. Massaro, "Development and Experimentation with Synthetic Visible Speech Behavioral Research Methods and Instrumentation," vol. 26, pp. 260-265, 1994.
- [16] R. Picard, "Affective Computing," 1997.
- [17] D. W. Massaro, "Illusions and Issues in Bimodal Speech Perception.," presented at Proceedings of Auditory Visual Speech Perception '98., Terrigal-Sydney Australia, December, 1998.
- [18] J. Lester, B. B. Stone, and G. Stelling, "Lifelike Pedagogical Agents for Mixed-Initiative Problem Solving in Constructivist Learning Environments.," *User Modeling and User-Adapted Interaction*, vol. 9, pp. 1-44, 1999.

- [19] J. Rickel and W. L. Johnson, "Animated agents for procedural training in virtual reality: Perception, cognition, and motor control.," *In Applied Artificial Intelligence*, vol. 13, pp. 343--382, 1999.
- [20] N. Person, A. Graesser, L. Bautista, E. M. and, and TRG, "Evaluating student learning gains in two versions of AutoTutor," *Artificial Intelligence in Education: AI-ED in the Wired and Wireless Future*, J. Moore, C. Redfield, and W. Johnson, Eds., pp. 286--293, 2001.
- [21] C. Lee and S. Narayanan, "Toward detecting emotions in spoken dialogs," *IEEE transaction on speech and audio processing*, vol. 13, 2005.
- [22] F. Dellaert, T. Polzin, and A. Waibel, "Recognizing Emotion in Speech," presented at Proceedings of the ICSLP, 1996.
- [23] C. Lee, S. Narayanan, and R. Pieraccini, "Classifying Emotions in Human-Machine Spoken Dialogs," presented at Proc. of International Conference on Multimedia and Expo, Lausanne, Switzerland, August 2002.
- [24] A. Paeschke and W. F. Sendlmeier, "Prosodic Characteristics of Emotional Speech: Measurements of Fundamental Frequency Movements," presented at Proceedings of the ISCA-Workshop on Speech and Emotion, 2000.
- [25] R. Tato, R. Santos, R. Kompe, and J. M. Pardo, "Emotional Space Improves Emotion Recognition," presented at Proc. Of ICSLP-2002, Denver, Colorado, September 2002.
- [26] S. Yacoub, S. Simske, X. Lin, and J. Burns, "Recognition of Emotions in Interactive Voice Response Systems," presented at The Eurospeech 2003, 8th European Conference on Speech Communication and Technology, Geneva, Switzerland, 1-4 September 2003.
- [27] F. Yu, E. Chang, Y. Q. Xu, and H. Y. Shum, "Emotion Detection From Speech To Enrich Multimedia Content," presented at the Second IEEE Pacific-Rim Conference on Multimedia, Beijing, China, October 24-26, 2001.
- [28] M. M. Louwerse and H. H. Mitchell, "Towards a taxonomy of a set of discourse markers in dialog: A theoretical and computational linguistic account.," *Discourse Processes*, pp. 199-239, 2003.
- [29] R. Duda, P. Hart, and D. Stork, *Pattern Classification*, 2nd ed. New York: Wiley, 2001.
- [30] I. H. Witten and E. Frank, *Data Mining: Practical machine learning tools and techniques*, 2nd ed. San Francisco: Morgan Kaufmann, 2005.
- [31] C. E. Williams and K. N. Stevens, "Emotions and speech: Some acoustical correlates," *JASA*, vol. 52, pp. 1238-1250, 1972.
- [32] R. Banse and K. R. Scherer, "Acoustic profiles in vocal emotion expression," *J. Personality and Social Psychology*, vol. 70, pp. 614-636, 1996.
- [33] S. Mozziconacci, "The expression of emotion considered in the framework of an intonational model," *Proc. ISCA Wrksp. Speech and Emotion*, pp. 45-52, 2000.
- [34] P. Boersma and D. Weenink, "Praat: doing phonetics by computer," Version 4.4.16 ed, 2006.
- [35] S. Kettebekov, M. Yeasin, and R. Sharma, "Prosody-based Audio Visual co-analysis for co-verbal gesture recognition," *IEEE transaction on Multimedia*, vol. 7, pp. 234-242, 2005.
- [36] M. Pantic and L. J. M. Rothkrantz, "Toward an affect-sensitive multimodal human-computer interaction.," *Proceedings of the IEEE*, vol. 91, pp. 1370 - 1390, Sept. 2003.