© 05-4506832 pustaka.upsi.edu.my Perpustakaan Tuanku Bainun Kampus Sultan Abdul Jelil Shah PustakaTBainun DEVELOPMENT OF AN AUTOMATIC ATTITUDE RECOGNITION SYSTEM: A MULTIMODAL ANALYSIS OF VIDEO **BLOGS**

Noor Alhusna Madzlan

Thesis submitted for the Degree of Doctor in Philosophy School of Linguistics, Speech & Communication Sciences

🕓 05-4506832 🔇 pustaka.upsi.edu.my f Perpustakaan Tuanku Bainun 🕥 PustakaTBainun 👘 ptbupsi

University of Dublin Trinity College

Under the Supervision and Direction of Asst. Prof. Breffni O'Rourke and Prof. Nick Campbell









Summary

Communicative content in human communication involves expressivity of socio-affective states. Research in Linguistics, Social Signal Processing and Affective Computing in particular, highlights the importance of affect, emotion and attitudes as sources of information for communicative content. Attitudes, considered as socio-affective states of speakers, are conveyed through a multitude of signals during communication. Understanding the expression of attitudes of speakers is essential for establishing successful communication. Taking the empirical approach to studying attitude expressions, the main objective of this research is to contribute to the development of an automatic attitude classification system through a fusion of multimodal signals expressed by speakers in video blogs. The present study describes a new communicative genre of self-expression through social media: video blogging, which provides opportunities for interlocutors to disseminate information through a myriad of multimodal characteristics. This study describes main features of this novel communication medium and focuses attention to its possible exploitation as a rich source of information for human communication. The dissertation describes manual annotation of attitude expressions from the vlog corpus, multimodal feature analysis and processes for development of an automatic attitude annotation system. An ontology of attitude annotation scheme for speech in video blogs is elaborated and five attitude labels are derived. Prosodic and visual feature extraction procedures are explained in detail. Discussion on processes of developing an automatic attitude classification model includes analysis of automatic prediction of attitude labels using prosodic and visual features through machine-learning methods. This study also elaborates detailed analysis of individual feature contributions and their predictive power to the classification task.

🕓 05-4506832 🜍 pustaka.upsi.edu.my 📑 Perpustakaan Tuanku Bainun 💟 PustakaTBainun 🚺 ptbupsi

Perpustakaan Tuanku Bainun Kampus Sultan Abdul Jalil Shał



ptbups

Related Publications

[1] Noor Alhusna Madzlan, Jingguang Han, Francesca Bonin and Nick Campbell.
Towards automatic recognition of attitudes: Prosodic analysis of video blogs.
Speech Prosody. Page 91-94, 2014

[2] Noor Alhusna Madzlan, Jingguang Han, Francesca Bonin and Nick Campbell.
Automatic recognition of attitudes in video blogs - Prosodic and visual feature analysis.
Fifteenth Annual Conference of the International Speech Communication Association, INTERSPEECH. Page 1826 - 1830, 2014

pustaka.upsi.edu.my Perpustakaan Tuanku Bainun Kampus Sultan Abdul Jalil Shah

[3] Noor Alhusna Madzlan, Yuyun Huang and Nick Campbell. Automatic classification and prediction of attitudes: Audio-visual analysis of video blogs. *Speech and Computer: 17th International Conference, SPECOM.* Springer. Volume 9319, Page 96 - 104, 2015

[4] Noor Alhusna Madzlan, Justine Reverdy, Francesca Bonin, Loredana Sundberg Cerrato and Nick Campbell. Annotation and Multimodal Perception of Attitudes: A Study on Video Blogs. *Third European Symposium of Multimodal Communication (MMSYM)*. Page 50-54, 2015



PustakaTBainun

05-4506832 😨 pustaka.upsi.edu.my 👔 Perpustakaan Tuanku Bainun Kampus Sultan Abdul Jalil Shah



Contents

1		Intro	duction	1	1
		1.1	Human	Communication	2
		1.2	Attitud	e as a Communicative Function	4
		1.3	Multim	odal Communication	6
			1.3.1	Multimodalities as Signals in Communication	6
		1.4	Multim	odal Affective Systems	7
			1.4.1	Use of Machine Learning Techniques	8
	05-4	1.5 ⁵⁰⁶⁸³² 1.6	Concep Motiva	pustaka.upsi.edu.my	10 tbupsi 11
		1.7	Statem	ent of Problem	12
		1.8	Researc	ch Objectives	12
		1.9	Detaile	d Outline	13
					10
2	2	State	of the	Art	15
		2.1	Comm	anicative Content	15
			2.1.1	Attitude - Definition	16
			2.1.2	Affective Attitudes	17
			2.1.3	Prosodic Attitudes	19
	,	2.2	Multim	odal Expression	24
			2.2.1	Prosodic Signals	25
			2.2.2	Visual and Facial Signals	28
			2.2.3	Fusion of Multimodal Expressions	31
			2.2.4	Multimodalities in Vlogs	33
	05-4	203832	Affecti	ye Recognition Systems roustakaan Tuanku Bainun Kampus Sultan Abdul Jalii Shah	tb36i

CONTENTS

	C	$) \begin{array}{c} 05.450683\\ 2.3.1 \end{array}$	Annotation Annotation Perpustakaan Tuanku Bainun Kampus Sultan Abdul Jalil Shah	PustakaTBainun	37 ptbupsi
		2.3.2	Machine Learning Techniques		40
	2.4	Concep	otual Underpinning		42
3	Vloa	Annots	ation Scheme		45
5	3 1	Introdu	ction		45
	2.1	The VI	cuon		45
	3.2			* * * * * * * *	40
		3.2.1	why vlogs?		46
		3.2.2	Vlog Characteristics		47
		3.2.3	Ethical Considerations	• • • • • • • •	50
		3.2.4	Speaker selection		52
		3.2.5	Video Selection	• • • • • • • • •	54
		3.2.6	Video Preparation		56
	3.3	Annota	tion and Segmentation		59
		3.3.1	Attitude Annotation Scheme		59
	C	3.3.268	Annotation Procedure ^{ny}	PustakaTBainun	61 ptbupsi
		3.3.3	Attitude Segmentation		63
	3.4	Validity	y of Attitudes		67
		3.4.1	Motivation		67
		3.4.2	Experimental Setup		68
		3.4.3	Results		72
		3.4.4	Discussion		79
	3.5	Conclu	sion		82
4	Mult	timodal	Feature Contribution		85
-	41	Introdu	ction		85
	1.2	Footuro	Extraction		02
	4.2	4.0.1			00
		4.2.1			ØØ
		4.2.2	Visual Features		90
	4.3	Feature	Selection		95
	C) 4:3 450683	Statistical Methods for Feature Selection	PustakaTBainun	960 ptbupsi

CONTENTS

C	05 4 4 2 06	68 Prominent Prosodic Features Rerpustakaan Tuanku Bainun Kampus Sultan Abdul Jalil Shah Pustaka TBainun	8 ^{si}
		4.4.1 Role of Pitch	1
		4.4.2 Role of Voice Quality	3
	4.5	Prominent Visual Features	7
		4.5.1 Role of Jaw	8
		4.5.2 Role of Eyebrows	8
	4.6	Fusion of Multimodal Features	2
		4.6.1 Multimodal Perception	2
		4.6.2 Multimodalities in Attitude Classification	4
	4.7	Discussion	7
	4.8	Conclusion	8
5	A	amotic Attitude Classification	-
Э	Auto	Introduction 12	1
	5.1		1
	5.2	Conceptual Applications	1
\bigcirc	5.4	Classification technique using SVMs Sultan Abdul Jali Shah	25
	5.4	Experimentation	4
		5.4.1 Experiment 1	5
		5.4.2 Experiment 2	9
		5.4.3 Experiment 3	4
	5.5	Discussion	9
	5.6	Conclusion	2
6	Con	nclusion 14	5
	6.1	Research Contributions	5
	6.2	Applications	7
	6.3	Research Constraints	8
	6.4	Future Work	0
A	17 12	IT *.4 . P T/* 3	
A	Full	List of videos 15	3
B	Füll	list of attitude segments f Perpustakaan Tuanku Bainun Kampus Sultan Abdul Jalil Shah	9 ^{si}

xii C	O5-4506832 Og pustaka.upsi.edu.my Perpustakaan Tuanku Bainun Guide to Creating a Vlog	CONTENTS PustakaTBainun ptbupsi 181
D	Annotation with Wavesurfer	185
E	Segmentation with Windows Live Movie Maker	189
F	Research Ethics Committee Approval Letter	195
G	Full list of 67 facial landmarks in AAM	197
H	Process of Visual Feature Extraction with AAM	201





O5-4506832 Sustaka.upsi.edu.my Perpustakaan Tuanku Bainun Kampus Sultan Abdul Jalil Shah

PustakaTBainun ptbupsi







O 5-4506832 pustaka.upsi.edu.my Perpustakaan Tuanku Bainun Kampus Sultan Abdul Jalil Shah PustakaTBainun



List of Tables

1.1	Atittude Components	4
2.1	19 affective attitudes	18
2.2	Attitudes in Past Literature	23
2.3	Modalities used to indicate types of information	24
2.4	Applying Prosody in Recognition Systems	27
2.5	Categories for Lip Tracking using ASM	29
2.6	Classification rate per modality from Kessous' work [92]	32
05-4506832	Example of MUMIN coding labels kaan Tuanku Bainun Pustaka Bainun	39 5005
3.1	Subset list of Video Information	54
3.2	Summary of Videos per Speaker	55
3.3	Standard A10-based Annotation Scheme	59
3.4	Five Attitude Labels	60
3.5	Segments by Attitude Category	64
3.6	Segments across Speaker	64
3.7	Subset List of Attitude Segments	66
3.8	Percentages of Selection for Each Attitude	74
4.1	Extracted Prosodic Features	87
4.2	Examples of Facial Landmark Labels	93
4.3	Prediction performance of Prosodic features	101
4.4	PC values of prosodic features	105
4.5	Predictive performance of Visual features	107
4.6 05-4506832	PC values of visual features . Perpustakaan Tuanku Bainun Pustaka TBainun	109 bupsi

LIST OF TABLES

5.1	Concepts for Supervised Learning Perpustakaan Tuanku Bainun Kampus Sultan Abdul Jalil Shah	122 ^{ptbups}
5.2	Total instances for each attitude label	125
5.3	Mean values for each attitude category with standard deviation in brackets .	127
5.4	Results for the different feature sets	128
5.5	Total instances for each attitude label	130
5.6	Prosodic Mean values for each attitude category with standard deviation in	
	brackets	131
5.7	Performance of the Trained Classifier	133
5.8	Total instances for each attitude label	134
5.9	Performance of Speaker-Independent Classification Task	135
5.10	Precision and Recall per Attitude Class	136
5.11	Confusion Matrix of Attitudes	137
5.12	Result of Visual Prediction	138
5.13	Summary of Experiments	138
A.1	List of videos 05-4506832 pustaka.upsi.edu.my f Perpustakaan Tuanku Bainun Kampus Sultan Abdul Jalil Shah PustakaTBainun	168 ptbupsi
B .1	List of Attitude Segments	180
G.1	List of Visual Features	199







05-

O5-4506832 Spustaka.upsi.edu.my Perpustakaan Tuanku Bainun Kampus Sultan Abdul Jalil Shah



List of Figures

1.1	Relationship between Affect and Emotion	3
1.2	Attitude as a Communicative Function	5
1.3	Flowchart for Supervised Machine Learning	9
2.1	Conceptual Framework	43
3.1	Characteristics of Video Blogs	48
3.2	Related Video Bloggers	53
3.3	Details of Audio File Format	57
3.4 ⁰	Details of Video File Format ^{Campus} Sultan Abdul Jalil Shah	tbupsi 58
3.5	Process of Annotation	62
3.6	Description of the A10 and N5 attitude categories	68
3.7	Participant Information	69
3.8	Sections in Online Survey	70
3.9	Example of Attitude Choices in Survey	71
3.10	Example of the 7-point Likert Scale of Certainty	72
3.11	Frequency of Occurrence for Attitude Selection	74
3.12	Barplots for Each Attitude's Certainty Rate	76
3.13	Barplots for Each "Other" Attitude's Certainty Rate	78
4.1	Acoustics measurements in the speech spectrum [65]	89
4.2	Facial Landmarks tracked by FaceSDK	92
4.3	Process of Facial Feature Extraction	94
4.4	Feature Selection in Supervised Machine Learning	95
4.5 06832	Average distribution of prosodic features pustaka.upsi.edu.my Rampus Sultan Abdul Jalil Shah	99 itbupsi

xvi		LIST OF FIGU	RES
(S) 05-4.6	4506832 pustaka.upsi.edu.my Perpustakaan luanku Bainun Kampus Sultan Abdul Jalil Shah Boxplot indicating distribution of Pitch	PustakaTBainun	ptbupsi 102
4.7	Boxplot indicating Voice Quality distribution		104
4.8	Scatterplot showing PCs for prosodic features		106
4.9	Scatterplot showing PCs for visual features		110
4.10	Distribution of Eyebrows		111
4.11	Agreed answers per modality		113
4.12	Feature selection using Decision Tree		115
5.1	Stages of Experimentation		124
5.2	Facial Tracking with FaceSDK AAM		132
C.1	Showing a vlog script		181
C.2	Related equipment needed		182
C.3	Filming a vlog		182
C.4	Editing raw film footage		183
C.5	Uploading video on YouTube		183
C.6 5-	Publishing video on YouTube Perpustakaan Tuanku Bainun Kampus Sultan Abdul Jali Shah	PustakaTBainun	184 ^{ptbupsi}
D.1	Open related wav file		185
D.2	Configure settings		185
D.3	Insert attitude label		186
D.4	Select related attitude label		186
D.5	Save transcription in directory		187
E.1	Open MP4 video file from directory		189
E.2	Select MP4 video file in relevant directory		190
E.3	Edit on menu bar		190
E.4	Trim tool		190
E.5	Mark start time		191
E.6	Mark end time		191
E.7	Save trimmed video		192
E.8	Save movie		192
E.9	Save to directory 4506832 pustaka.upsi.edu.my F Perpustakaan Tuanku Bainun Kampus Sultan Abdul Jalil Shah	PustakaTBainun	193 ptbupsi

H.1	Open AAM software	201
C H.2	Select Video 06832 yustaka.upsi.edu.my f Perpustakaan Toanku Bainun Pustaka TBainun Pustaka TBainun Pustaka TBainun	201 ptbupsi
H.3	Facial Tracker	202
H.4	Output folder	202
H.5	Rename Output folder	202
H.6	Structure of data	203

05-4506832 🚱 pustaka.upsi.edu.my F Perpustakaan Tuanku Bainun Kampus Sultan Abdul Jalil Shah







05-4506832 Spustaka.upsi.edu.my Perpustakaan Tuanku Bainun Kampus Sultan Abdul Jalil Shah

PustakaTBainun ptbupsi

List of Abbreviations

Abbreviation	Meaning	Page
SVM	Support Vector Machine	9
RF	Random Forest	9
VLOG	Video Blog	10
AES	Affective-Epistemic States	18
F0	Fundamental Frequency	19
HNR	Harmonicity-to-Noise ratio	19
B1	Bandwidth of the First Formant	25
H1	Amplitude of the First Harmonic	26
Al	Amplitude of the First Formant	26
H1-H2	First harmonic relative to the Second Harmonic	26
H1-A3	First harmonic relative to the Third Formant	26
NAO	Normalised Amplitude Ouotient	26
ASM	Active Shape Model	28
FAC	Facial Action Coding system	29
AAM	Active Appearance Model	29
DCTustaka.upsi.d	Discreet Cosine Transform	ustaka RO in
HMM	Hidden Markov Model	30
ANOVA	Analysis of Variance	31
HCI	Human Computer Interaction	32
wMEI	Weighted Motion Energy Images	33
GP	Gaussian Process	36
GMM	Gaussian Mixture Model	40
LDA	Linear Discriminant Analysis	41
Vlogger	Video Blogger	17
Hz	Hertz	80
dB	Decibel	80
SEC	Second	80
FMFAN	Mean value of the fundamental frequency	80
FMIN	Minimum value of the fundamental frequency	80
FMAX	Maximum value of the fundamental frequency	80
FPCT	Percentage of shape of the pitch contour	80
FVCD	Percentage of the vibration of the vocal folds	00
PMFAN	Mean value of power/ intensity	00
PMIN	Minimum value of the nower/ intensity	00
PMAY	Maximum value of the power/ intensity	00
PPCT	Intensity movement	80
DN	Duration	80
PCA	Principal Component Applying	00
PC	Principal Components	90
PC1	First Principal Component	90
DC2	Second Principal Component	97
r CZ	Second Principal Component	97
	Automotio Speech Dece	97
ASK	Automatic Speech Recognition	100
KDF	Radial Basis Function	114





Chapter 1

Introduction

This dissertation highlights affective expressions as a source of information content in human communication. The study explores the dynamics of attitudinal expressions conveyed through several types of non-verbal signals. Attitudinal expressions between communicators must be understood and displayed appropriately to ensure successful message transfer during the communication process. This research further extends understanding of attitudinal aspects of expression, not only through the exploration of human perception of attitudes, but also through the development of automatic attitude recognition which will be of great utility for automatic understanding of user perception and information retrieval.

This chapter briefly discusses human communication, its different signalling modalities for communicating content of various types, as well as the role of attitudinal states in the communicative setting. Further discussion about the meaning of attitudes and their relation, or non-relation, to affect and emotions is also explored. Several definitions and concepts are introduced to give a clearer overall view of the methodology involved in the study. The final part of this chapter explains the research motivation and objectives.





1.1 Human Communication

Perpustakaan Tuanku Bainun Kampus Sultan Abdul Jalil Shah

Research in domains including social psychology, social signal processing and affective computing has increasingly focused on understanding the dynamics of human communication. Humans communicate with each other in a number of different modalities and the interpretation of their utterances may vary from a person to the other. Generally, humans communicate with each other to share information. Human communication involves a series of processes and can take place between two or more interlocutors. This sharing of information is transferred through signals.

Allwood [1] elaborates on the purpose of human communication, which is dissemination of information. Allwood mentions several types of information involved in the communication process:

- 1. physiological states (fatigue and hunger)
- 2. character, identity, personality (being timid, aggressive)

3. affective-epistemic attitudes (showing joy, friendliness, surprise)

4. factual content giving information about beliefs, assumptions about facts

5. communication management (feedback, turn-taking)

Among the information streams used in communication, one interesting type of information is the expression of affective-epistemic attitudes. This type of information is useful for understanding affective states and emotions of the people involved in the communicative setting. Emotions, affect and attitudes seem to be similar concepts but there are clear distinctions between them. Damasio [2] distinguishes between emotion and feelings, stating that feelings refer to the inner, cognitive experience of an emotion while emotion is the observable response of these feelings. The frequently cited five basic emotions are sadness, anger, fear, surprise, disgust and happiness [3]. There are certainly other emotional ways of representing emotional states than the simply categorical, including valence based representations [4] but such concepts are not sufficient to describe the complexities of people's affective states [5]. A person could express a mixture of sometimes contradictory feelings simultaneously.



job, she expresses first the feeling of happiness and shifts her expression to sadness. Hence, human emotion, although easily shown, are often complex and confusing to understand or interpret.

Affect, on the other hand, is considered a general term for the inward feelings of the human experience. This concept is a broader representation of feelings where emotion contributes to a large part of the overall definition. While Zanna and Rempel [6] refer to affect as "any thoughts that are infused with strong, weak or no emotion at all", Shouse [7] claims that affect has no relation with feelings and emotion. The term affect is sometimes used interchangeably with emotion. It is believed that there are differences between the two concepts. Emotion refers to the display of feelings while affect is a non-conscious experience. Affect is a pre-personal experience of the speaker that unconsciously affects the consequent feelings and actions of the speaker. Hence, affect is viewed as a broad, abstract concept of the humanly cognitive experience, while emotion refers to the inner feelings of the person. Figure 1.1 summarizes the relationship between affect and emotion.



Figure 1.1: Relationship between Affect and Emotion

Affect is a general and broader representation of a person's feelings while emotion is the displayed response of affect, and thus is incorporated into the affective state of a person.

Recent empirical studies, particularly in the field of Affective Computing, use affect to refer to affective-epistemic states of humans when interpreted by machines. This use of affect is a broad term as a point of reference to studies on emotion, attitude and behavioural states of humans displayed in the communicative setting. This study supports the concept of attitude as representation of speakers' affective-epistemic characteristics. Attitude is viewed differently from concepts of affect and emotion. The following section details attitude as a

concept for study. pustaka.upsi.edu.my

Perpustakaan Tuanku Bainun Kampus Sultan Abdul Jalil Sha



1.2 Attitude as a Communicative Function Sustaka TBainun

The previous section explains the differences between affect and emotion. The concept of attitude is related to affect and emotion, but attitude has distinctive traits. Zanna and Rempel [6] describe attitude as states that may be expressed as strong emotions or may be identified solely from the way an individual behaves with an object. Research on the psychology of attitude finds limited agreement on precise definitions and characteristics of attitude. Psychologists relate attitude with beliefs, opinions, habits and values [8]. Oskamp [8] explains attitude as having three main components, as shown in Table 1.1:

Components	Description
Cognitive	Ideas and beliefs of the agent towards the object
Affective (emotional)	Feelings and emotions the agent has towards the object
Behavioural	The agent's actions towards the object

Table 1.1: Atittude Components

Fishbein and Ajzen [9] suggest that attitude is mainly associated with the affective compoputaka.upsi.edu.ny Kampus Sultan Abdul Jalit Shah Putaka Balmun nent. This is believed to be true as these components are independent and separate entities but are still interrelated. The affective component of attitude is an object of interest to several fields of research including psychology, economics, and marketing. However, some researchers use the term attitude loosely and often interchangeably with affect or emotion. In this study, these concepts, although similar, are treated independently. Attitude is of particular interest for this study as it represents observable traits of speaker's cognitive experience. Attitude refers to actions, outer representation of feelings, while emotion refers to inner feelings that are difficult to evaluate. Attitude is believed to be a pragmatic concept for the understanding of affect and emotions [10]. This interpretation of emotional states involves the crucial aspect of intended, voluntary and controlled actions. Auberge [10] explains communicative functions in Figure 1.2.

Perpustakaan Tuanku Bainun Kampus Sultan Abdul Jalil Shah



4



Figure 1.2: Attitude as a Communicative Function

Figure 1.2 shows attitudinal functions as a component for communication. The scale indicates emotional functions to be at the far left of the control scale while linguistic functions are on the far right. Attitudinal function is in the middle of the control scale. This means that attitudes, unlike emotions, are not vague and involuntarily expressed. Rather values of speaker intention are expressed in a conscious, controlled and voluntary manner.

The definition of attitude may or may not have a strong emotional component but always has an evaluative component. In fact, attitude and emotion can be seen as a continuum, a degree of emotional involvement in attitudes. Emotion, conversely, need not have a strong attitudinal - i.e. evaluative - component. This is because emotions are not necessarily directed towards any particular content, that is, you can be happy or sad (for example) without being able to identify why. Additionally, emotions can be internal and less susceptible to conscious control. Whereas attitudes are held more consciously, often deliberately expressed (in language), and subject to control in the sense of reflective evaluation.

Emotion and attitude are conceptually distinct, but interact with one another in reality. Attitudes always have an element of evaluation. Oskamp [8] has identified three components of attitude; cognitive, affective/emotional and behavioural. In this thesis I am concerned with the overt expression of attitude, which I take to be evidence of the first two, internal, components, i.e., the cognitive and affective/emotional states of the speaker. Under one understanding of the term attitude, attitudes can be entirely internal, comprising the cognitive and affective states, but never overtly expressed. For the purposes of this research, the focus is on attitudes as overtly displayed. Therefore, when I use the term attitude, it will frequently 05-4506832 pustaka.upsi.edu.my refer to the overt expression of attitude.

1.3 05- Multimodal Communication Abdul Jalil Shah

Attitude is expressed in the content of communication through several different channels or modalities. Multimodal communication is defined as "co-activation, sharing and coconstruction of information simultaneously and subsequently through several modes of perception and production" [11]. This concept adds to the concepts of communication with the sharing of information as its main function. However, multimodal communication involves several modes of information sharing using simultaneous sensory channels. These include sight, hearing, touch, smell and taste. Multimodal communication involves two main processes; firstly multimodal integration (perceiving information based on several modalities) and multimodal distribution (producing information using multimodalities). This is similar to the traditional communicative process, where both communicative agents perceive information and produce feedback. The difference is that there is emphasis on simultaneous sensory modalities.

Multimodal communication has become an important research area as humans perceive and produce communicative expressions through several modalities. A combination of speech upsi visuals and gesture facilitates communicators to meet their communicative goal of information sharing. As mentioned in Section 1.1, the contents of information include affectiveepistemic attitudes. The simultaneous use of multimodal signals facilitates expression of attitudes. Due to the complexities of affective attitudes, the use of more than one sensory modality is beneficial to better understand and interpret attitudes. For example, a person can perceive a friend's expression of friendliness through rising tone of voice, display of a smile and the wave of a hand. Simultaneous use of speech, facial expressions and gestures facilitates the communicators' message transmission.

1.3.1 Multimodalities as Signals in Communication

Numerous studies suggest the relation between display of affective states and non-verbal gestures [12] [13] [14] [15]. Vinciarelli and Valente [16] refer to non-verbal communication as the transmission of a message through non-verbal behavioural cues, such as facial expressions, vocalisations, gesture and posture. Non-verbal communication relays speakers' inner feelings, whether intentionally or unintentionally. For instance, a mother outwardly thupsi

6

9 PustakaTBainun

expresses approval when her child voluntarily makes the bed by establishing eye contact with the child, praising her with a rising tone and giving a bright smile. In another scenario, the mother tries not to show disapproval when the child throws a tantrum by speaking with a level or falling tone and controlling her facial expressions from showing contempt. This example shows how people communicate feelings and intentions through the dynamics of several and simultaneous non-verbal signals.

With the emergence of the body of knowledge of Social Signal Processing, the role of this area is to firstly provide physical quantification and synthesis of non-verbal signals through which the affective behaviour of humans are expressed, and, secondly, to implement non-verbal signals in conversational agents [12]. This area is especially interesting as multimodal traits of humans are quantified through development of recognisers and synthesisers. This study aims to evaluate and measure communicative contents, in particular attitudinal representation of speakers. The outcome can be helpful in facilitating better understanding of relationships in human-human and human-computer interaction.

1.4 Multimodal Affective Systems

One method of quantifying attitudinal states of speakers through multimodal signals is by applying machine-learning methods. Research in the fields of affective computing, speech technology and human-computer interaction employs machine-learning techniques. There are numerous studies conducted in emotion recognition based on analysis of sentiment [17] [18], prosody [19], facial features [15], gestures [20] and posture [21]. Multimodality (using combination of signals) is also studied in great detail in emotion recognition. However, when a distinct separation of terminology between emotions and attitudes is made, there is little research conducted on automatic attitude recognition. Although there are some notable studies on prosodic attitudes [22] [23][24], research in automatic classification of attitudes using multimodal signals is still scarce.

In this century, there is growing interest in social media as a rich source of human expression. People from different areas of study, including sociolinguistics, psychology, affective computing and human-computer interaction conduct studies using data from social

media. Social psychologists, for example, investigate people's communicative purposes in D5-4506832 pustaka.upsi.edu.my Perpustakaan Tuanku Bainun Kampus Sultan Abdul Jalil Shah



online blogs [25] [26] [27]. Researchers are also interestered in the social activity of microblogging on Facebook [28] [29]. Other empirical studies involve research in online search engines [30] [31], Twitter [32] [33] and YouTube [34] [35] [36]. Through conceptual and empirical research, it is evident that social media offers a rich source of social discovery and scientific application. Similarly, open access to instances of human behaviour, in particular, attitude expressions, from social media users shows a dynamic social representation worthy of exploration and understanding.

1.4.1 Use of Machine Learning Techniques

Machine learning is used in the area of Artificial Intelligence to enable computers to learn without being explicitly programmed. The learning takes place through observations of new data, hence the outcome of this implementation is a system that enables automatic classification of the learned data. Essentially, there are several methods of machine learning – supervised, semi-supervised and unsupervised learning. Supervised learning involves the task of training a dataset that has a group of features and object labels. The machine is con-05-4506832 Supervised and unsupervised features and object labels. The machine is constructed to predict and identify the label of an object given the set of features. Unsupervised learning however functions without the identification of labels. Instead, labels are derived based on the training of the dataset given by the features provided.

Classification in supervised learning is used to predict categorical labels given a set of observations. Figure 1.3 [37] illustrates the work flow of building a classification model using supervised learning technique:



Figure 1.3: Flowchart for Supervised Machine Learning

Figure 1.3 generally describes the processes involved in developing a predictive model using ⁵⁻⁴⁵⁰⁶⁸³² pustaka upst.edu.my a supervised machine learning method. Following this training process, a predictive model is derived and will produce expected labels.

There are several learning algorithms developed for machine learning classification tasks. One popular method of implementation is the use of statistical learning algorithms [38]. This study adopts this method by using Support Vector Machine (SVM) and Random Forest (RF) classifiers, as will be explained in Chapter 5.



1.5 05-4506832 pustaka.upsi.edu.my Conceptual Definitions

This section addresses conceptual definitions used in this thesis. The list describes definitions for the following terminologies:

- 1. Annotation: Annotation is a metadata that is associated to another data by commenting and noting. In this study, annotation involves the act of assigning attitude labels on relevant parts of the video blogs. Annotation of attitude states creates another set of data comprising of only video segments that are annotated with an attitude label.
- Attitude Expression: Attitude refers to pragmatic interpretations of emotion [10]. The expression of attitudes are socio-affective states of speakers expressed in voluntary and controlled settings.
- 3. **Machine-Learning**: Machine-learning is an area of study in Artificial Intelligence where automatic algorithms are constructed to allow learning and prediction of computers from the given data.

4. **Prosody**: Prosody refers to the suprasegmentals in the voice. The use of prosody in this study covers aspects of fundamental frequency and pitch contours, intensity of the voice, voice quality and duration of speech segments.

- 5. **Visual**: Visual refers to facial observations from videos used in this study. This reference to visual features includes movement of the eyes and other facial expressions of the speaker's face.
- 6. Vlog: Vlog is an abbreviated term for video blog. This abbreviation is used in this research to refer to YouTube videos of speakers where they share stories of daily life and events.









10

05-4506832 pustaka.upsi.edu.my **1.6 Motivation**

There has been much work in the fields of Social Signal Processing and Affective Computing in developing recognisers for affective states of speakers through multimodal signals. Human beings can detect differences in affective expressions using data gathered through their eyes and ears. Machines can also do the same using information recorded by cameras and microphones. So by physically capturing and then analysing affective states of speakers, researchers are able to understand human behaviour through concrete displays of affective expressions. Treating attitude as dissimilar to affect and emotion, the development of an attitude recognition system through the means of multimodal signals is an interesting step forward towards understanding human behaviour. Apart from that, development of automatic recognisers has two main purposes:

First Automatic information retrieval and categorisation. Recognition systems are used to index and retrieve information from metadata or other forms of media. Collection of information obtained from these systems enables other systems to implement this source of information for affective modelling. Perpustakaan Tuanku Bainun Kampus Sultan Abdul Jalil Shah

Second Applications for human-machine interaction. This system could be applied to artificial communicative agents such as social robots for technological advancement. Communicative content can be interpreted by an automatic system and used to inform the behaviour to robots or avatars. This is useful to then develop socially intelligent robots that understand and respond appropriately to humans.

This study aims to contribute to the first purpose of developing automatic recognition systems. Communicative content, in particular, affective information is useful to be embedded into systems for information retrieval and categorisation. This creates a source for people to retrieve information on attitude expressions during speech.

Attitude expressions are especially prevalent in social media where people find comfort and openness in expressing themselves better than in face-to-face communication. One advantage of social media is the ease of public access. This allows researchers to investigate human behaviours, particularly affective expressions from a dynamic source. YouTube, for example, allows public access for people to share videos about their daily life. This rich human behaviour. Treating these videos as representations of people's attitude expressions makes the collection of this corpus interesting for development of an attitude recognition system that could be applicable to other areas of study.

1.7 Statement of Problem

To the best of my knowledge, little research has so far been conducted in the area of automatic attitude recognition, particularly in exploiting social media corpus such as YouTube video blogs to understand multimodal expressions of humans. This study introduces the concept of attitudes as an evaluative representation of human social experience through spontaneous talk in vlog speech. Following that, internalising attitude states in a computer system for artificial intelligence also contributes to this study's major work by applying multimodal signals into the recognition system.

1.8 Research Objectives

v pustaka.upsi.edu.my

The main objective of this research is to develop an automatic attitude recognition system for information retrieval and categorisation. To achieve this, the following objectives are defined:

Objective 1 Collection of a corpus that represents attitude expressions. Expressions of attitudes are observed and investigated from social media. The use of social media is essential as they provide dynamic and rich multimodal signals that include speech, facial expressions and gestures to indicate speakers' attitude expressions. This study annotates and segments five attitude expressions from vlogs in YouTube.

Objective 2 Investigation of the use of multimodal features as contributors to developing a reliable attitude recognition system. Relevant prosodic and visual features are extracted and analysed. Subsequently, specific prosodic and visual features are selected and examined to understand their contribution towards the classification model.

Objective 3 Development of a reliable automatic attitude recognition system. This research highlights the development of an attitude recognition system that can predict different atti-

attitude expressions through multimodal signals, namely prosodic and visual features. Supervised machine learning is conducted using Support Vector Machines (SVM) and Random Forests (RF).

1.9 Detailed Outline

This thesis is divided into six main chapters which are listed below:

Chapter 1 Discussion of general concepts of human communication and multimodalities in affective expression of attitudes. This chapter also outlines research motivation, statement of problems and three main research objectives for the present study. Conceptual definitions used in this thesis are also listed in this chapter.

Chapter 2 Elaboration of discussion and criticisms of theories of attitudes, multimodalities and use of social media for understanding human communication. This chapter also outlines some of the recent works in the study of automatic recognition through multimodalities using machine learning techniques.

Chapter 3 Discussion of a collection of the vlog corpus. Vlogs are collected from an online video sharing website, YouTube. An ontology of attitude annotation scheme is introduced in this chapter. The processes of collection, annotation and segmentation of attitude states of speakers are described in detail.

Chapter 4 Discussion of the processes involved in developing an automatic attitude classification system. This process involves steps for multimodal feature extraction, from prosodic and visual means. The second part of the chapter discusses processes of feature selection. This section highlights prominent prosodic and visual features that provide greatest contribution to the classification system in recognising different attitudinal states of speakers.

Chapter 5 Supervised machine-learning techniques and results on building a reliable attitude classification system. A total of three experiments are reported:

a) development of an attitude classification system through prosodic modality

b) use of prosodic and visual modalities to develop an automatic attitude classification model

c) improved use of multimodal feature sets for improving the classifier

Chapter 6 Statement of general conclusions, research limitations, suggestions for improve-05-4506832 Perpustaka.upsi.edu.my ment and future direction of this Study. Perpustakaan Tuanku Bainun Kampus Sultan Abdul Jalil Shah



Perpustakaan Tuanku Bainun Kampus Sultan Abdul Jalil Shah



ptbups

Chapter 2

State of the Art

This dissertation focuses on developing an attitude recognition system for the purpose of classification and retrieval of human attitudes from video recordings. Prior to more elaborate discussion, this chapter presents the state-of-the-art from related literature addressing the notion of attitude, its relevance to multimodal and novel forms of media, and recent applications in recognition interfaces. The present work adopts some of the conceptual frameworks and methods used from past literature to address the research goals. This chapter begins with a description of related literature concerning attitude definitions and their relation to communicative content.

2.1 Communicative Content

As we have noted, Allwood [1] sees the exchange of information between communicators. We will now explore the five types of communicative information that he identifies (repeated here for convenience):

- 1. Physiological states
- 2. Character, Identity and Personality
- 3. Affective-Epistemic Attitudes

pustaka.upsi.edu.my

4. Factual Contents

() 05-4506832

5. Communication / Feedback Management

Perpustakaan Tuanku Bainun Kampus Sultan Abdul Jalil Sha





Physiological states are expressed in the communicative setting, and these include fatigues and hunger. The character of a person is also communicated in speech, examples include the aggressiveness and openness of an individual. Affective-epistemic attitudes refer to the affective or emotive state of the speaker: for instance friendliness, joy and surprise. Examples of factual content include information about beliefs, theories and assumptions, while communication management consists of turn-taking, sequences of feedback and topic change [1]. Another aspect of communicative content that is relevant to the present work in this thesis is Affective-Epistemic Attitudes. This content is essential for the communicative process during information transfer. The next section describes relevant literature relating to attitudes.

2.1.1 Attitude - Definition

Communicative content during communication involves several types of information. One crucial source of information, as mentioned by Allwood [1] is the affective-epistemic attitudinal states of communicators. To better understand the communicative process as well as the modes of interaction, it is essential to first understand the role that attitudes have [1]. In the literature, attitudes are interchangeably related to affect and emotions. However, this thesis identifies attitudes as a distinctive term and this is described through the definition and components of attitudes.

Malhotra [39] states distinctions between affect and attitude, which agrees with Ajzen that attitude refers to summary evaluations of an object or behaviour [9] [40]. Bargh and Chartrand explain attitude as being judgments that are outcomes of spontaneous and unconscious effort [41]. They claim that attitude can be automatically activated without any prior goals of judgment [41]. This perspective, however, is not supported by Auberge [10] who views attitude as driven by intention, voluntary and controlled action. Auberge describes a continuum (see Figure 1.2) of controlled expression of affect where emotion is regarded as an inward and involuntary cognitive experience. Attitude is pragmatic representations of emotions, voluntary and controlled expressions of the cognitive state while linguistic function is the linguistic means of affective expression. Hence, attitude stems from a pragmatic puscha upsi edu my for Perputakaan Tuanku Bainun (See Figure 1) of controlled action, and not from a spontaneous and unconscious effort.