

**USING TEXT MINING ALGORITHM TO DETECT GENDER DECEPTION
BASED ON MALAYSIAN CHATROOM LINGO**

BY

**ASSOCIATE PROFESSOR DIANNE L.M. CHEONG
NUR ATIQAH SIA BINTI ABDULLAH@SIA SZE YIENG**

OCTOBER 2006



UNIT PENYELIDIKAN PROTRAD
FAKULTI TEKNOLOGI MAKLUMAT DAN SAINS KUANTITATIF
UNIVERSITI TEKNOLOGI MARA
40450 SHAH ALAM
E-mail: protrad@tmsk.itm.edu.my
Homepage: <http://protrad.tmsk.itm.edu.my>

Tarikh : 29 September 2005
Surat kami: 100-PT/RU/3/6-(43/05)

Prof. Madya Dianne L. M. Cheong
Pensyarah (Ketua Projek Penyelidikan)
Fakulti Teknologi Maklumat dan Sains Kuantitatif
UiTM, Shah Alam

Puan,

CADANGAN PENYELIDIKAN:
USING TEXT MINING ALGORITHM TO DETECT GENDER
DECEPTION BASED ON MALAYSIAN CHATROOM LINGO

Dengan segala hormatnya perkara di atas adalah dirujuk.

Sukacita dimaklumkan bahawa ProTRAD pada 21 September 2005 telah membuat keputusan:

- i. Bersetuju meluluskan cadangan penyelidikan yang dikemukakan oleh Puan dan Nur Atiqah Sia Abdullah.
- ii. Tempoh projek penyelidikan ini ialah **12 bulan**, iaitu mulai 1 November 2005 hingga 31 Oktober 2006.
- iii. Kos yang diluluskan ialah sebanyak RM 8,000.00 sahaja.
- iv. Penggunaan geran yang diluluskan hanya akan diproses setelah perjanjian ditandatangani.
- v. Semua pembelian peralatan yang kosnya melebihi RM 500.00 perlu menggunakan Pesanan Jabatan Universiti Teknologi MARA (LO). Pihak Puan juga dikehendaki mematuhi peraturan penerimaan peralatan.
- vi. Kertas kerja boleh dibentangkan dalam seminar setelah **75% draf awal Laporan Akhir** projek dihantar ke Unit Penyelidikan ProTRAD untuk semakan. Walau bagaimanapun, permohonan kepada Institut Penyelidikan, Pembangunan dan Pengkomersilan (IRDC) untuk tujuan pembentangan perlu dibuat terlebih dahulu.
- vii. Pihak Puan dikehendaki mengemukakan Laporan Kemajuan Projek Penyelidikan bagi tempoh 4 bulan pertama, iaitu sehingga 28 Februari 2006, 4 bulan berikutnya iaitu sehingga 30 Jun 2006 dan 4 bulan berikutnya iaitu sehingga 31 Oktober 2006. Laporan akhir

ABSTRACT

E-mail is used for communication between strangers and friends. It can be a fantasy playground for identity experimentations where players take on an imaginary persona and interact with each other in the virtual world. In communication, knowing the identity of those whom you communicate is essential for understanding and evaluating an interaction. However, the presentation of self in the virtual world is often a conscious and deliberate endeavour. Therefore, gender deception is difficult and risky and it can be abandoned at will. Inference can be made both from writing style and from clues hidden in the posting data. A text-mining algorithm was designed to detect gender deception based on gender-preferential features at the word or clause level of Malaysian e-mail users. Based on this designed text algorithm, a prototype in Visual Basic is developed. The prototype was tested with 16 documents; each consists of 5 e-mails exchanges of respective individuals. Out tests have shown that the prototype is at 81.3% of accuracy level. This is consistent with a human reader of the documents. The tested prototype will be a tool to assist interest parties such as the Criminology and Forensic Department, e-mail users and virtual communities to successfully identify gender deception.

Keywords: *gender detection, gender of e-mail author, text mining algorithm to detect gender, program to detect gender, gender deception.*

TABLE OF CONTENTS

	PAGE
LETTER OF APPROVAL	ii
LETTER OF TRANSMITTAL	iv
PROJECT TEAM MEMBERS	v
ACKNOWLEDGEMENT	vi
ABSTRACT	vii
TABLE OF CONTENTS	viii
LIST OF TABLES	x
LIST OF FIGURES	xi
 CHAPTER 1 INTRODUCTION	 1
1.1 Background of the Study	1
1.2 Rationale of the Study	2
1.3 Statement of the Problem	4
1.4 Objectives of the Study	5
1.5 Research Questions	6
1.6 Significance of the Study	6
1.7 Limitation of the Study	7
1.8 Definitions of Terms	8
1.9 Summary	8
 CHAPTER 2 LITERATURE REVIEW	 10
2.1 Related Studies on Text Mining Algorithm	10
2.1.1 TDM	11
2.1.2 GATE	13
2.1.3 D2K	14
2.1.4 YALE	14
2.1.5 CEM	15
2.2 Related Studies on Gender Deception on the Internet	15
2.2.1 Online systems	17
2.2.2 Instant communication	17
2.3 Theoretical Framework of the Study	18
2.3.1 E-mail	18
2.3.2 Chat rooms	19
2.3.3 Instant Messaging	20
2.4 Conceptual Framework of the Study	21
2.4.1 E-mail Text body	21
2.4.2 Text mining algorithm	26
2.4.3 Visual Basic program	26
2.4.4 Gender of author of e-mail	28
2.5 Summary	29

CHAPTER 3 METHODOLOGY	30
3.1 Content Analysis	30
3.2 Content Analysis by Features	33
3.3 Text Mining Algorithm Design	34
3.3.1 Text Body	35
3.3.2 Feature Set Extractor	35
3.3.3 Item Feature Set	36
3.3.4 Evaluation	36
3.3.5 Inference	36
3.4 Summary	37
CHAPTER 4 PROGRAM DESIGN AND CODING	38
4.1 Flash Screen	38
4.2 Main Form	38
4.3 Statistic Form	53
4.4 Result Form	54
4.5 Help Form	54
CHAPTER 5 IMPLEMENTATION AND TESTING	55
5.1 Prototype Manual	55
5.2 Testing	60
CHAPTER 6 FINDINGS AND CONCLUSION	77
6.1 Discussion on Findings	77
6.2 Recommendation to the Prototype	79
6.3 Future Research Works	79
6.4 Conclusion	80
REFERENCES	81
APPENDIX A E-mail Documents	83