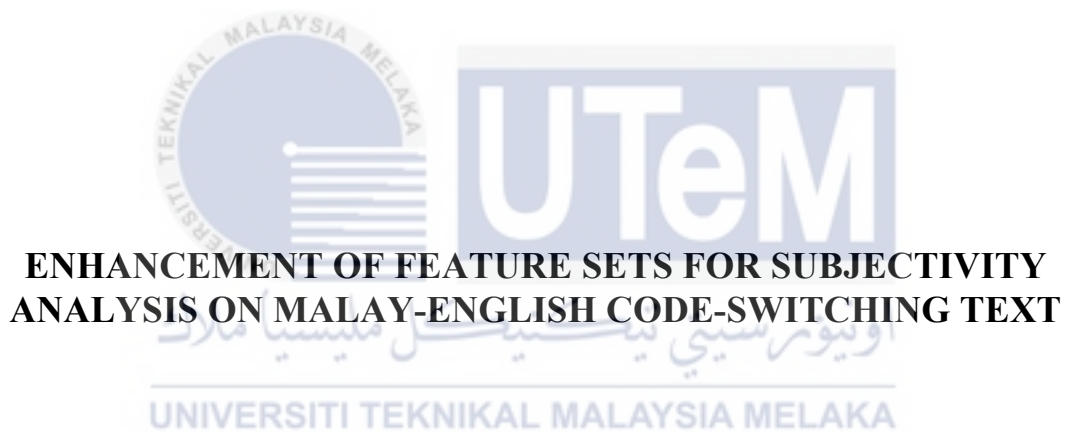




Faculty of Information and Communication Technology



**ENHANCEMENT OF FEATURE SETS FOR SUBJECTIVITY
ANALYSIS ON MALAY-ENGLISH CODE-SWITCHING TEXT**

Emaliana Binti Kasmuri

Doctor of Philosophy

2023

**ENHANCEMENT OF FEATURE SETS FOR SUBJECTIVITY ANALYSIS ON
MALAY-ENGLISH CODE-SWITCHING TEXT**

EMALIANA BINTI KASMURI

**A thesis submitted
in fulfilment of the requirements for the degree of Doctor of Philosophy**



Faculty of Information and Communication Technology

UNIVERSITI TEKNIKAL MALAYSIA MELAKA

2023

DECLARATION

I declare that this thesis entitled “Enhancement Of Feature Sets For Subjectivity Analysis On Malay-English Code-Switching Text” is the result of my own research except as cited in the references. The thesis has not been accepted for any degree and is not concurrently submitted in candidature of any other degree.



Signature :

Name : EMALIANA BINTI KASMURI

Date : 11 January 2023

UNIVERSITI TEKNIKAL MALAYSIA MELAKA.....

APPROVAL

I hereby declare that I have read this thesis and in my opinion, this thesis is sufficient in terms of scope and quality for the award of Doctor Philosophy.

Signature : *lieza*

Supervisor Name : DR. HALIZAH BINTI BASIRON

Date : 11 January 2023



اونيورسيتي تيكنيكل مليسيا ملاك
UNIVERSITI TEKNIKAL MALAYSIA MELAKA

DEDICATION

To my beloved husband, father, mother and sister.

Thank you for all your love, support and patience.



ABSTRACT

A code-switching sentence is a sentence that is constructed using two or more languages. It is a norm for a multi-lingual speaker to use code-switching sentences to share objective and subjective textual information on public platforms such as blogs and social media. Classifying a voluminous code-switching text into subjective and objective classes has posed a new challenge to the current solution of subjectivity analysis. The current solution has limited its design to process only monolingual text. Therefore, the presence of subjective code-switching text is ignored by the current solution. The ignorant limits the capability of the current solution to generate an accurate result of subjectivity analysis on code-switching text. Therefore, this research aims to find a set of solutions for subjectivity analysis on code-switching text. The research process begins by filling in the absence of the subjectivity code-switching corpus. A subjective Malay-English code-switching corpus was built. The corpus contains 35,067 Malay-English code-switching sentences that were harvested from Malay-English blog posts. Each sentence was annotated with either subjective or objective labels. The research process continues with designing the feature sets that represent the subjectivity of the Malay-English code-switching sentences from the corpus. The feature sets were enhanced from the subjective monolingual feature set, that was initially designed to represent subjectivity of English text. The initial subjective monolingual feature sets consist of pronoun, adjective, cardinal number, modal and adverb. The enhanced feature sets consist three feature sets which are embedded code-switching feature set, unified code-switching feature set and stylistic feature set. The embedded code-switching feature used the initial monolingual feature set for English and embeds the feature of Malay language in it. In the unified code-switching feature set, the extracted Malay and English features were unified using an adapted algorithm known as the Malay-English Unified POS. The algorithm predicts the type of each word in a code-switching sentence according to the language of the word. In the stylistic feature set, emoticons, interjections, signs of subjectivity such as exclamation marks and word with exaggerations of spelling were extracted to represent the subjectivity in the code-switching sentences. The effectiveness of the enhanced feature sets was evaluated using the Malay-English code-switching subjectivity corpus as the data set and two machine learning classifiers, which are Naïve-Bayes and Support Vector Machine. The 10-fold cross-validation classification technique was used on different settings of experiments and combinations of feature sets to obtain the performance of the enhanced feature sets. The performance from the combination of unified code-switching and stylistic feature sets has outperformed other feature sets. The combination has consistently performed at the accuracy of 59% using both machine learning classifiers. The consistent performance indicates the combined feature sets are the viable solution for subjectivity analysis on the Malay-English code-switching text.

**PENAMBAHBAIKAN SET FITUR UNTUK ANALISA SUBJEKTIVITI KE ATAS
TEKS PERALIHAN KOD MELAYU-INGGERIS**

ABSTRAK

Sepotong ayat peralihan kod ialah ayat yang dibina menggunakan dua atau lebih bahasa. Adalah menjadi satu kebiasaan bagi seorang yang fasih pelbagai bahasa menggunakan ayat-ayat peralihan kod berkongsi maklumat tekstual objektif dan subjektif di atas platform terbuka seperti blog dan media sosial. Pengelasan teks peralihan kod yang banyak kepada kelas-kelas subjektif dan objektif telah memberikan cabaran baru kepada penyelesaian sedia ada bagi analisa subjektiviti. Rekabentuk penyelesaian sedia ada telah menghadkan rekabentuk penyelesaian hanya untuk memproses teks bahasa mono. Oleh itu, kehadiran teks subjektif peralihan kod diabaikan oleh penyelesaian sedia ada. Pengabaian ini telah menghadkan keupayaan penyelesaian sedia ada untuk menjana keputusan yang tepat bagi analisa subjektiviti ke atas teks peralihan kod. Maka, matlamat kajian ini adalah untuk mendapatkan satu set penyelesaian bagi analisa subjektiviti ke atas teks peralihan kod. Proses kajian ini bermula dengan mengisi kelompongan korpus subjektiviti peralihan kod. Sebuah korpus subjektiviti peralihan kod Melayu-Inggeris telah dibangunkan. Korpus tersebut mengandungi 35,067 ayat-ayat peralihan kod Melayu-Inggeris yang dituai dari hantaran-hantaran blog. Setiap ayat dilabelkan dengan subjektif atau objektif. Proses kajian diteruskan dengan mereka bentuk set-set fitur yang mewakili ciri-ciri subjektif bagi ayat-ayat peralihan kod Melayu-Inggeris. Set-set fitur ini ditambahbaik daripada set fitur subjektif bahasa mono, yang mana pada asalnya mewakili subjektiviti teks Bahasa Inggeris. Set fitur awal subjektif bahasa mono mengandungi kata ganti nama, adjektif, nombor kardinal, modal dan adverba. Set-set fitur yang ditambahbaik mengandungi the set fitur iaitu set fitur peralihan kod tertanam, set fitur peralihan kod kesatuan dan set fitur gayaan. Set fitur peralihan kod tertanam menggunakan set fitur awal bahasa mono untuk Bahasa Inggeris dan ditanam fitur bahasa Melayu di dalamnya. Di dalam set fitur peralihan kod kesatuan, fitur-fitur bahasa Melayu dan Inggeris yang telah disatukan menggunakan algoritma yang telah diadaptasi yang dikenali sebagai POS Kesatuan Melayu-Inggeris. Algoritma ini meramal jenis setiap perkataan di dalam ayat peralihan kod berdasarkan kepada bahasa perkataan tersebut. Di dalam set fitur gayaan, emotikon, celahan, tanda subjektiviti seperti tanda seru dan perkataan dengan ejaaan yang berlebihan digunakan untuk mewakili ciri subjektiviti di dalam ayat-ayat peralihan kod. Keberkesanan penambahbaikan set-set fitur dinilai menggunakan korpus subjektiviti peralihan kod Melayu-Inggeris sebagai set data dan dua pengelas pembelajaran mesin iaitu Naïve-Bayes dan Support Vector Machine. Teknik klasifikasi bersilang 10-lipatan digunakan ke atas pelbagai tetapan dan kombinasi eksperimen bagi mengukur pencapaian ke atas penambahbaikan set-set fitur. Pencapaian daripada kombinasi set fitur kesatuan dan gayaan telah mendahului set-set fitur yang lain. Kombinasi tersebut telah mencapai ketepatan 59% secara konsisten menggunakan kedua-dua pengelas pembelajaran mesin. Pencapaian yang konsisten menunjukkan kombinasi set-set fitur adalah penyelesaian yang berdaya maju untuk analisa subjektiviti ke atas teks peralihan kod Melayu-Inggeris.

ACKNOWLEDGEMENTS

In the name of Allah, the Most Gracious and the Most Merciful

I thank Allah for granting me the patience, health, guidance and determination to complete this thesis successfully. This thesis exists because of the assistance, support and inspiration of many people. Firstly, I acknowledge both Universiti Teknikal Malaysia Melaka dan the Government of Malaysia, for supporting me and giving me the opportunity.

I'd want to extend my heartfelt gratefulness to Dr Halizah Basiron for her encouragement and assistance during my work on my thesis. She introduced me to a variety of online and library options. She addressed all of my inquiries and assisted me in narrowing down my search. I appreciate how she constructed her constructive and positive critique and how attentively she reviewed my work. I also wish to thank my second supervisor Dr Mohamed Ishak Desa, who retired during the completion of my studies, for being patient with me and my questions, and Dr Yogan Jayakumar for his comments on my work.

Lastly, for all my friends who support me at Fakulti Teknologi Maklumat dan Komunikasi, I gave them my thanks and pray that Allah will protect them from harm.

TABLE OF CONTENTS

	PAGE
DECLARATION	
APPROVAL	
DEDICATION	
ABSTRACT	i
ABSTRAK	ii
ACKNOWLEDGEMENTS	iii
TABLE OF CONTENTS	iv
LIST OF TABLES	viii
LIST OF FIGURES	xi
LIST OF APPENDICES	xiv
LIST OF ABBREVIATIONS	xv
LIST OF PUBLICATIONS	xvii
CHAPTER	
1. INTRODUCTION	1
1.0 Introduction	1
1.1 Research Background	1
1.2 Problem Statement	4
1.3 Research Question	6
1.4 Objectives	8
1.5 Scope of Study	9
1.6 Significance of Study	9
1.7 Research Contribution	11
1.8 Structure of Thesis	13
1.9 Summary	14
2. LITERATURE REVIEW	15
2.0 Introduction	15
2.1 Subjectivity Analysis	15
2.1.1 Types of Information in Textual Document	17
2.1.2 Subjective Information	18
2.2 Subjectivity Classification vs Polarity Classification	22

2.3	Subjectivity Classification	25
2.3.1	The Role of Subjective Lexicon	29
2.3.2	The Early Work of Subjective Classification	34
2.4	Subjectivity Classification in Various Genres of Documents	36
2.4.1	Subjectivity Classification for Formal Written Style Document	38
2.4.2	Subjectivity Classification for Casual Written Style Documents	49
2.4.3	Subjectivity Classification for Blogs	50
2.4.4	Subjectivity Classification for Web-Based Forum	55
2.4.5	Subjectivity Classification for Review	59
2.4.6	Subjectivity Classification for Tweets	61
2.5	Issues in Subjectivity Classification	64
2.6	Code-Switching Text Analysis	66
2.6.1	Definition of Code-Switching	66
2.7	Monolingual Sentences vs Code-Switching Sentences	68
2.7.1	Construction of Code-Switching Sentences	70
2.8	Attempt of Code-Switching Text Analytical Studies	74
2.9	Annotation of Code-Switching Corpus	77
2.10	Challenges in Text Analytical Studies on Code-Switching Document	81
2.11	Summary	83
3.	RESEARCH METHOD	86
3.0	Introduction	86
3.1	Research Design	86
3.2	Research Operational Framework	87
3.2.1	Phase 1: Preliminary Study	88
3.2.2	Phase 2: Code-Switching Subjectivity Corpus Construction	90
3.2.3	Phase 3: Code-Switching Subjectivity Enhanced Features Construction	91
3.2.4	Phase 4: Code-Switching Subjectivity Classification	93
3.2.4.1	The Naïve Bayes Classifier	93
3.2.4.2	The Support Vector Machines Classifier	96
3.2.4.3	Evaluation of Performance	98
3.2.5	Phase 5: Research Findings Documentation	102
3.3	Summary	102
4.	BUILDING A CODE-SWITCHING SUBJECTIVITY CORPUS	103
4.0	Introduction	103
4.1	Overview of the Code-Switching Corpus Development	104
4.1.1	Harvesting Blog Posts	107
4.1.2	Extracting Code-Switching Sentences	110

4.1.3	Result of Code-Switching Sentence Extraction	115
4.2	Instantiating the Annotation Theory	117
4.3	Selecting and Training The Annotators	120
4.4	Specifying Annotation Procedure	120
4.5	Choosing Annotation Tool	121
4.6	Choosing and Applying Evaluation Measure	122
4.7	Delivery and Maintaining the Corpus	124
4.8	Summary	125
5.	PROPOSED ENHANCED FEATURE SETS FOR SUBJECTIVITY CLASSIFICATION ON MALAY-ENGLISH CODE-SWITCHING TEXT	126
5.0	Introduction	126
5.1	Code-Switching Subjectivity Analysis Methodology	127
5.2	Text Pre-processing	128
5.2.1	Token Validator	130
5.2.2	Token Normalization	133
5.2.3	Token Part-of-Speech (POS) Tagger	134
5.3	Design of Feature set for Text Classification	135
5.4	Feature Design and Representation	137
5.4.1	Initial Feature Set	137
5.5	Enhanced Subjectivity Feature Set for Code-Switching Sentences	138
5.5.1	Embedded Code-Switching Feature Set	140
5.5.2	Unified Malay-English POS Feature Set	147
5.5.3	Stylistic Feature Set	151
5.6	Summary	153
6.	RESULT AND DISCUSSION	154
6.0	Introduction	154
6.1	Description of Dataset	155
6.2	Experiment Setup	157
6.3	Discussion of the Experimental Result	158
6.3.1	Baseline Feature sets Performance Results	159
6.3.2	Embedded Code-Switching Feature set Performance Results	161
6.3.3	Unified Code-Switching POS Feature set Performance Results	166
6.3.4	Stylistic Feature sets Performance Results	169
6.3.5	Combinations of Feature Sets Performance Results	171
6.4	Summary	176

7. CONCLUSION AND RECOMMENDATION OF FUTURE WORKS	177
7.0 Introduction	177
7.1 Contribution of the Study	178
7.1.1 The Malay-English Code-Switching Subjectivity Corpus	178
7.1.2 Embedded Code-Switching Feature Set	179
7.1.3 Unified Malay-English Code-Switching Feature Set	181
7.1.4 Stylistic Feature Set	183
7.1.5 Rule-Based Language Identification for Code-Switching Text	184
7.1.6 Other Contributions	184
7.2 Recommendation of Future Works	185
7.3 Summary	187
REFERENCES	188
APPENDICES	204



LIST OF TABLES

TABLE	TITLE	PAGE
1.1	Summary of Problem Statements	6
1.2	Mapping of Research Question to Problem Statement	7
1.3	Mapping of Objectives and Research Questions	9
2.1	Definition of private states dimension	16
2.2	List of commonly used adverbial linguistic items and part-of-speech (POS)	28
2.3	Summary of studies in subjectivity classification designed for English formal written document	44
2.4	Summary of studies in subjectivity classification design for non-English formal written document	48
2.5	Example of noise in casual written style document	49
2.6	Summary of studies in subjectivity classification for blog	54
2.7	Summary of subjectivity classification for web-based forum	58
2.8	Examples of feature used to represent the content-free features	60
2.9	Summary of studies in subjectivity classification for review	61
2.10	Summary of studies in subjectivity classification using tweets	64
2.11	Summary of subjectivity classification of different genres with impressive performance	65
2.12	Comparison of code-switching pattern	73
2.13	Summary of studies for text analysis on code-switching documents	76
2.14	Reliability value of sentiment analysis datasets	78

2.15	Comparison of reliability measurement	80
2.16	Comparison of monolingual text and code-switching text	82
3.1	Interpretation of confusion matrix	100
4.1	Basic statistics concerning the downloaded blog posts	108
4.2	Number of sentences based on total of words per sentence	109
4.3	Total number of sentences by type of sentence	115
4.4	Summary of Annotation Scheme for Malay-English Code-Switching Corpus	119
4.5	Kappa value interpretation	123
4.6	Samples of annotated sentences from MS-EN-CS corpus	124
5.1	Examples of valid tokens	130
5.2	Examples of invalid tokens	130
5.3	List of retained punctuations	133
5.4	Permissible repeated alphabet character and its position	134
5.5	Representation of feature set	136
5.6	Initial feature set for subjectivity analysis	138
5.7	Sample of results from English POS Tagger on Malay-English code-switching sentences	142
5.8	Sample of results from Malay POS Tagger on Malay-English code-switching sentences	143
5.9	Grouping of Malay and English Part-of-Speech (POS) Tags	145
5.10	Embedded code-switching subjectivity feature set	146
5.11	Representation of emoticon and its underlying meaning	152
5.12	Example of interjections and its meaning	153
6.1	Performances of initial feature set models	161
6.2	Performance of embedded feature sets	163

6.3	Performance of unified feature sets	167
6.4	Performance of stylistics feature sets	170
6.5	Performance results for subjectivity classification on the Malay-English code-switching text	175



LIST OF FIGURES

FIGURE	TITLE	PAGE
2.1	Relationship of subjectivity analysis and sentiment analysis	17
2.2	Example of (1) subjective sentence and (2) objective sentence	18
2.3	Granularity of textual document processes by the sentiment analysis system	21
2.4	Comparison of published articles between subjectivity analysis classification and polarity classification between 2007 and 2017	23
2.5	Process of subjectivity classification	26
2.6	Example of synset for happy WordNet Search 3.1	31
2.7	Genres of documents analyzed for the studies of subjectivity classification	38
2.8	Example of mixed language sentence	67
2.9	Example of inter-sentential code-switching	68
2.10	Example of monolingual sentences in (a) Malay and (b) English	69
2.11	Name of entity in foreign language	70
2.12	Example of code-switching text	70
2.13	Example of insertion pattern	71
2.14	Example of alternation pattern	72
2.15	Example of congruent lexicalization pattern marked in orange colored font	72
3.1	Research operational framework	89

3.2	Flow of activities in Phase 2	91
3.3	Example of projected feature vectors into a space	96
3.4	Example of hyperplane for two subspaces	97
3.5	Confusion matrix	99
4.1	Flow of creating the subjectivity corpus adapted	106
4.2	Distribution of sentences according to the number of words per sentences	110
4.3	The flow of process to extract the Malay-English code-switching sentences	111
4.4	The rule-based procedure to compute distribution of Malay-English words in a sentence	112
4.5	A procedure to determine Malay-English code-switching sentence	114
4.6	Distribution of extracted code-switching sentences	116
4.7	Distribution of annotated sentences	122
5.1	Workflow in code-switching subjectivity analysis	128
5.2	Procedure of text pre-processing	129
5.3	Procedure for token validator	132
5.4	Proposed feature set for code-switching subjectivity analysis	139
5.5	Visualization of enhanced code-switching feature set	140
5.6	Flow of subjectivity analysis for embedding Malay POS	141
5.7	The flow of subjectivity analysis for Malay-English code-switching sentences using unified POS feature set	148
5.8	Procedure for English-Spanish POS Unification from Solorio and Liu (2008b)	149

5.9	Procedure to generate unified code-switching POS Tag	150
6.1	Distribution of dataset for the experiment	156
6.2	Accuracy performances of baseline initial features models using different classifiers across multiple datasets	159
6.3	Result of accuracy performance for subjectivity classification on Malay-English code-switching text using embedded feature set	162
6.4	Comparison of embedded feature sets with baseline initials feature sets	164
6.5	Results of accuracy performance for subjectivity of classification on Malay-English code-switching text using the unified feature set	166
6.6	Comparison of accuracy performance for unified feature set with baseline initial feature sets	168
6.7	Results of accuracy performance for subjectivity of classification on Malay-English code-switching text using the unified feature set	170
6.8	Comparison of accuracy performance for stylistic model with baseline initial feature sets	171
6.9	Comparison of averaged accuracy performances for combined feature sets using the Naïve-Bayes classifier	173
6.10	Comparison of averaged accuracy performances for combined feature sets using the SVM classifier	174

LIST OF APPENDICES

APPENDIX	TITLE	PAGE
A	List of English Stop Words	204
B	List of Malay Stop Words	211
C	List of Translated Modals	219
D	Profile of Annotators	220
E	Profile of Language Experts	225



LIST OF ABBREVIATIONS

AMI	-	Augmented Multi-party Interaction
AUROC	-	Area Under the Curve
BERT	-	Bidirectional Encoder Representations from Transformers
BiLSTM	-	Bidirectional Long Short-Term Memory
BOW	-	Bag-of-Words
CNN	-	Convolutional Neural Network
CRF	-	Conditional Random Field
D-BiLSTM	-	Double Bidirectional Encoder Representations from Transformers
DSE	-	Direct Subjective Expression
EM	-	Expectation Maximization
EN	-	English
ESE	-	Expressive Subjective Element
EWGA	-	Entropy Weighted Genetic Algorithm
GloVe	-	Global Vector
HTML	-	Hypertext Markup Language
IDE	-	Integrated Development Environment
IG	-	Information Gain
LAMP	-	Logistic Average Misclassification Percentage
LDA	-	Latent Dirichlet Allocation

MAP	-	Mean average precision
MPQA	-	Multi Perspective Question Answer
MS-EN	-	Malay-English
MS-EN-CS	-	Malay-English Code-Switching
NLTK	-	Natural Language Toolkit
OQ	-	Opinion Queries
POS	-	Part-of-Speech
PMI	-	Pointwise Mutual Information
RQ	-	Research Question
SLM	-	Subjective Language Model
SVM	-	Support Vector Machine
SWSD	-	Subjective Word Sense Disambiguation
TREC	-	Text Retrieval Conference
TF-IDF	-	Term Frequency-Inverse Document Frequency
UKM	-	Universiti Kebangsaan Malaysia
URL	-	Uniform Resource Locator
U.S	-	United States
VSM	-	Vector Space Model
WMI	-	Web-based Mutual Information

LIST OF PUBLICATIONS

Kasmuri, E. and Basiron, H., 2020. Segregation of Code-Switching Sentences using Rule-Based Technique. *International Journal of Advances in Soft Computing and its Applications*, 12 (1), 49–64.

Kasmuri, E. and Basiron, H., 2019. Building a Malay-English code-switching subjectivity corpus for sentiment analysis. *International Journal of Advances in Soft Computing and its Applications*, 11 (1), 112–130.

Kasmuri, E. and Basiron, H., 2017. Subjectivity analysis in opinion mining - A systematic literature review. *International Journal of Advances in Soft Computing and its Applications*, 9 (3), 132–159.



CHAPTER 1

INTRODUCTION

1.0 Introduction

This chapter gives an overview of the research work carried out for subjectivity analysis on mixed-language documents known as code-switching documents. It begins with a description of the research background. The research problems in the subjective analysis of code-switching documents were highlighted. Research questions were drawn from the problems and objectives were outlined to fulfil the aim of the research. The output of the research is listed. The organization of this thesis is described before the end of the chapter. This chapter is closed with a summary.

1.1 Research Background

Information is easily accessible today as compared to many decades ago. Users create information on various platforms, such as personal websites, blogs and social media. The information can be factual or non-factual, or a mixture of both. Factual information describes truths, details, or particulars about a subject matter related to a text (Banea et al., 2014). For example, facts about a book include the title, the number of pages, the book's author, the synopsis and the book's price. Non-factual information contains evaluative or affective information from some aspect of the described subject matter, such as opinion about the storyline from the perspective of a reader, for example, good or bad, or comments from a reader on the design of the book cover, for example, pleasant or unpleasant. Factual

information is known as objective information, whereas non-factual information is known as subjective information (Banea et al., 2014).

Subjective and objective information has become a piece of important information in decision-making. Both parts of information are used for different purposes. For example, subjectivity analysis was used in an automated question-and-answer system to interpret the type of questions asked, either seeking answers for personal opinions or experiences or verifying information (Li et al., 2008). Subjectivity analysis is also used to analyse the political stance during political campaigns aiding the party in planning its campaign strategy (Jiang and Argamon, 2008). Interestingly, subjectivity analysis was also used in automatic article classification to identify and extract sports betting information that will increase the chances of winning bets in the sports gambling (Hajj et al., 2019). Therefore, this information needs to be segregated into subjective and objective information.

The task of segregating the subjective and the objective information has become burdensome to the decision maker when a voluminous of information is presented to them. The segregation of subjective and objective information becomes more challenging because this kind of information is usually unstructured, where formats are not specified before the creation of the information (Chaturvedi et al., 2018). The information creators are not bound to any specific sentence structure, choice of words or flow of information when it was created. In addition, some bilingual or multilingual speakers have the habit of mixing different languages in the text when expressing their thoughts or opinions.

The habit of mixing different languages in a text has increased the challenge to separate subjective information from objective information (Çetinoğlu et al., 2016). Mixing different languages in a single sentence is known as the code-switching (Muysken, 2000). Code-switching is a common scenario for most bilingual or multilingual speakers and writers. The term code-mixing was also used to describe the mixing of different languages for a