# Natural Language Processing Based Features for Sarcasm Detection: An Investigation Using Bilingual Social Media Texts

Mohd Suhairi Md Suhaimin[1,2], Mohd Hanafi Ahmad Hijazi[1], Rayner Alfred[1,3] and Frans Coenen[4]

[1]Faculty of Computing and Informatics, Universiti Malaysia Sabah, Malaysia

[3]Centre of Excellence in Semantic Agent, Universiti Malaysia Sabah, Malaysia

[2]Human Resource Management Division, Ministry of Higher Education, Malaysia

[4]Department of Computer Science, University of Liverpool, United Kingdom

mohd.suhairi@mohe.gov.my, {hanafi, ralfred}@ums.edu.my, coenen@liverpool.ac.uk

*Abstract*—**The presence of sarcasm in text can hamper the performance of sentiment analysis. The challenge is to detect the existence of sarcasm in texts. This challenge is compounded when bilingual texts are considered, for example using Malay social media data. In this paper a feature extraction process is proposed to detect sarcasm using bilingual texts; more specifically public comments on economic related posts on Facebook. Four categories of feature that can be extracted using natural language processing are considered; lexical, pragmatic, prosodic and syntactic. We also investigated the use of idiosyncratic feature to capture the peculiar and odd comments found in a text. To determine the effectiveness of the proposed process, a non-linear Support Vector Machine was used to classify texts, in terms of the identified features, according to whether they included sarcastic content or not. The results obtained demonstrate that a combination of syntactic, pragmatic and prosodic features produced the best performance with an F-measure score of 0.852.**

*Keywords—sarcasm detection; feature extraction; bilingual feature; sentiment analysis*

## I. INTRODUCTION

The presence of sarcastic content has been hindering the performance of current Sentiment Analysis (SA) systems for some time. The issue is that in the presence of sarcasm the polarity of user expressions is reversed from its literal meaning resulting in misclassification, thus affecting the performance of SA systems [1]. The detection of sarcasm remains a research challenge although work has been directed at extraction of features indicative of sarcasm, the creation of corpuses featuring sarcasm, and mechanisms for detecting sarcasm such as classification mechanisms [2].

The issue of sarcasm detection is made more complex in the case of social media texts written in more than one language. The misspelled words, shortened word forms and stylistic text coupled with the use of dual language is commonplace in Malaysian social media texts, where it is not unusual to mix Malay and English. For example: "*Gaji minimum 900 dulu pun bukan semua laksanakan. Inikan pulak 1000. Just in ur dream lah.*" and "*tsk tsk tsk. xlme lg yuran nek la tu. good job.*". Detecting sarcastic contents in such bilingual situation adds an extra level of complication to the challenge of sarcasm detection. The crucial part is to extract the features into usefulness to represent the text to perform the sarcasm detection.

In this paper, we proposed a process to extract features that are indicative of the presence of sarcasm based on Natural Language Processing (NLP) with respect to bilingual social media texts, namely Malay and English languages. NLP has shown the ability to identify the presence of sarcasm and resolve the challenge [3] by means of recognizing and extracting a specific feature to simplify the complex meaning of texts [4]. The proposed process comprises two phases. In the first phase, we extract the lexical, pragmatic and prosodic features, in the order that they appear in a given bilingual text. In the second phase we translate the bilingual dataset to English and extract further prosodic features along with syntactic and idiosyncratic features. A Bag of Words (BOW) representation was used as the fundamental document representation because this is the simplest, and therefore the most commonly encountered even though the approach ignores grammar, word order, sentence structure and usually punctuation, thus ignoring the linguistic structure [5]. On its own, it is clearly not expressive enough to provide information regarding the presence of sarcasm, hence the overlay of syntactic and lexical features as proposed in this paper.

The main contributions of this paper are: (i) a process for extracting features indicative of the presence of sarcasm in bilingual social media texts and (ii) a set of NLP based feature categories appropriate for sarcasm detection in the context of Malay social media data. The rest of this paper is organized as follows. Section 2 presents the related work. Section 3 describes the proposed extraction process and Section 4 specifics the experiments conducted. Section 5 discusses the result and analysis. Section 6 concludes this paper and provides the future work.

## II. RELATED WORK

Sarcasm is a subtype of verbal irony; simplistically it is something which is said that should be interpreted as having the opposite meaning to its literal meaning [6]. Along with hyperbole, jocularity, rhetorical questions and understatements [7]; sarcasm is intended to convey a mixture of both obvious and subtle interpersonal meaning [8]. In terms of usage and interpretation, the nature may vary according to language and geographical area [9]. Feature extraction in the presence of sarcasm is particular significance in the context of SA. In this respect a number of different categories of feature that may be extracted using NLP techniques, have been proposed; examples include: lexical, pragmatic, prosodic and syntactic features.

Lexical features are the most common form of feature extracted using NLP. These features are used to represent the sentence's content using tokenized words in the form of n-grams, that can uncover information for processing [3]. Three types of n-gram are commonly used: 1-grams (unigrams), 2-grams (bigrams) and 3-grams (trigrams). Pragmatic features are intended to emphasis the meaning of the content of sentences that may include sarcasm [10]. Emoticons, so called "heavy" punctuation, hashtag (#) and repeated words are examples of pragmatic features. Prosodic features involve different pitches, loudness, timing and tempos in writing [11]. Interjections are an example of prosodic features. Syntactic features play an important role in providing information concerning the syntactic structure of documents. A common syntactic feature is the Part of Speech (POS) tag that may be associated with words in a document. All these categories of feature are likely to play a role with respect to the identification of features indicative of sarcasm.

Idiosyncrasy is a mode of behavior or way of thought peculiar to an individual[1]. In the linguistic study of metaphorical language, an idiosyncrasy is an isolated metaphor rarely used in common conversation and yet intended to bring meaning to the overall message. Idiosyncratic features have been used as non-systematic metaphorical expressions as part of "message delivery" during conversation. Examples include 'head of cabbage', 'foot of the mountain', and 'leg of a table' [12]. The point of metaphor is that the actual meaning should not be interpreted as the literal meaning (as in the case of sarcasm). Idiosyncratic features may be considered to be an indicator of the presence of sarcasm.

Some research has been conducted on sarcasm detection for non-English language data, examples can be found in [10],

[11], [12], [13] and [14]. Early work on sarcasm detection in the context of the Indonesian language, and using lexical features combined with syntactic and prosodic features on Twitter data, can be found in [13]. The English SentiWordNet[2] was translated into Indonesian using Google Translate[3]; only words existent in the translated SentiWordNet were considered. Lexical features in the form of unigrams were then extracted from the texts. The features extracted include negation, word context, affix, number of interjections and question words. Experiments, using the Support Vector Machine (SVM) model, recorded a best accuracy of 54.1% for sarcasm detection using a negation and interjection feature combination, outperforming the Naïve Bayes (NB) and Maximum Entropy (MaxEnt) classifiers. In other work on sarcasm detection using a Dutch Twitter dataset, only lexical feature were used (unigrams, bigrams and trigrams) with a frequency count of more than three [14]. The collected dataset was based on the Twitter hashtag *#sarcasme* (sarcasm in Dutch), thus signaling the presence of sarcasm. The hashtag *#sarcasme* was used as sarcastic marker that alerted the reader to the fact that the text was sarcastic. The features were then weighted using the chi-squared metric to select meaningful features. The weighted rankings demonstrated that exclamation occurred most frequent in Tweets with the *#sarcasme* label. Experiments using Balanced Winnow classifier [15], recorded a best Area Under Curve (AUC) of 0.79 in a balanced distribution and 0.75 in an imbalanced distribution.

In addition to lexical features, syntactic features have been employed to detect sarcasm using Czech Twitter data [16]. The lexical features used include various n-gram and frequency patterns while the syntactic features used were POS tags, namely: nouns, verbs, and adjectives, as well as the ratios of nouns to adjectives and adverbs. Lexical features with frequency patterns were found to produce the best result. A best F-measure (Fm) of 0.582 was recorded using a SVM classifier, outperforming the MaxEnt classifier in the reported evaluation. More recent work directed at detecting irony, a more general form of sarcasm, using Greek political Tweets has proposed various types of feature, including lexical, semantic, prosodic and pragmatic features, for sarcasm detection [17]. The features were grouped into four categories: lexical, spoken (stylistic writing), rarity (frequency of rare words), emoticon and meaning (measures of ambiguity using Greek WordNet[4] synsets). The groups were then ranked using Information Gain. The rarity feature group was ranked top followed by lexical. The authors employed supervised and semi-supervised classifications using the proposed features. A best Fm score of 0.79 was produced using supervised classification and Functional Trees. Functional Trees maximized the univariate and multivariate from of decision trees with a linear function as used in constructive induction learning. Although much work have been directed at non-English language feature detection, as far as the authors are aware there is no work directed at bilingual data.

---

[1] https://en.oxforddictionaries.com/definition/idiosyncrasy

[2] http://sentiwordnet.isti.cnr.it/
[3] https://translate.google.com/
[4] http://compling.hss.ntu.edu.sg/omw/

Based on the above, the five most promising categories of feature, extracted using NLP, for sarcasm detection in various languages were: lexical, pragmatic, prosodic, syntactic and idiosyncratic. The previous work, as reported above, also indicates that supervised classification was good sarcasm detection mechanism.

## III. THE PROPOSED FEATURE EXTRACTION PROCESS

We present in this section the features that can be extracted using NLP, that we intend to adopt for sarcasm detection in the context of bilingual data. In this paper, we considered both the original bilingual corpus and its translation in English for feature extraction. The process of feature extraction consists of two main steps: (i) extraction of the lexical, pragmatic and Malay prosodic features from the original bilingual corpus (Sub-section A), and (ii) translation of the bilingual corpus to English and extraction of English prosodic features, along with syntactic and idiosyncratic features (Sub-section B). The process returns sets of extracted features as summarized in Table I.

TABLE I.    TYPES OF FEATURE EXTRACTED

| Feature | Types | Corpus |
|---------|-------|--------|
| 1. Lexical | Unigram | Original bilingual |
| 2. Pragmatic | Punctuation marks, Hashtag | Original bilingual |
| 3. Prosodic | Interjection | Original bilingual and English |
| 4. Syntactic | Part of Speech | English |
| 5. Idiosyncratic | Idiosyncratic | English |

### A. Extraction from Bilingual Corpus

This section describes the process of feature extraction from the original bilingual corpus. Prior to the extraction of features, the corpus was preprocessed. The processes involved three steps as presented in Sub-sections 1 to 3 below.

#### 1) *Preprocessing of the Corpus and Lexical Features Extraction:*

Preprocessing a given corpus involved tokenization and spellchecking. Tokenization breaks the corpus into words and symbols such as punctuation and hashtag (#). Social media data contain highly noisy text which includes spelling errors, non standard words, stylistic words, short form words and repetitions. Classification accuracy is affected as the presence of noise increases [18]. Phrases such as "*xlme lg*", "*x lame lg*", "*tak lama lg*" or "*tak lama lagi*" are examples of texts that contain short form words, stylistic words and spelling errors; the correct standard phrase is "*tidak lama lagi*". The presence of the above cause "dispersion", where features that should be considered to be the same feature are treated as different features, result in poor performance when creating training data with which to build a classifier [19]. This we used Malay and English dictionaries to correct misspelled words. The preprocessing of the corpus also involved stopword removal to eliminate meaningless words. We used Malay and English stopword lists[5,6] for our bilingual data. Lexical feature were then extracted from the corpus in the form of n-grams. Single character such as 'n', 't', and 'b' were omitted. We then lowercased all the tokens.

#### 2) *Pragmatic Features Extraction:*

Punctuation marks were considered to be pragmatic features, instead of sentence segmentators, because of their potential to indicate sarcasm [20]. Heavy punctuation, for example high occurrences of various punctuation marks, is often an indicator of the presence of sarcasm in text. The punctuation marks considered in this work were: question marks (?), exclamation marks (!) and quotation marks ("" and ''). In addition hashtags (#) also considered. The length of sequences of punctuation marks was reduced to a maximum of three characters to avoid dispersion.

#### 3) *Prosodic (Malay) Features Extraction:*

In the third step a Malay list of interjections[7] was employed. It should be noted that interjections differs according to language, for example "*ooi*", "*puii*" and "*weii*" are only found in Malay. A total of 43 Malay interjections were identified and used. Malay prosodic features were thus extracted from the preprocessed corpus using this Malay interjection list.

### B. Extraction of Features from English Translated Corpus

This section describes the process of feature extraction from the English translated corpus.

#### 1) *Corpus Translation to English:*

The preprocessed corpus, generated as described in Sub-section A.1, was translated into English using Google Translate[8]. Although the resulting translations were by no means perfect, they were judged to produce translations that were sufficiently accurate to support further analysis, better than the translations using Moses or Bing [21]. The employment of deep learning recently boost the performance of Google Translate [22]. The translation preserved some Malay words such as names, locations and abbreviations; examples were 'Idris', 'Barisan Nasional', 'Sabah' and 'PTPTN'.

#### 2) *Prosodic (English) Features Extraction and Combination:*

The English interjection's list was obtained from publically available sources[9]. A total of 66 English interjections were employed. The interjections that already listed in the Malay interjection as described in Sub-section A.3 were removed to avoid redundant features. English prosodic features were then extracted from the translated corpus using the English interjection's list. The English prosodic features were then combined with the Malay prosodic feature from Section A.3.

#### 3) *Syntactic Features Extraction:*

---

[5] http://nlp.cs.nyu.edu/GMA_files/resources/malay.stoplist
[6] http://www.nltk.org/book/ch02.html
[7] https://ms.wikipedia.org/wiki/Kata_seru
[8] https://translate.google.com/
[9] http://grammar.yourdictionary.com/parts-of-speech/interjections/list-of-interjections.html

We choose four groups of POS: NOUN, VERB, ADJECTIVE and ADVERB based on Japerson's Theory for ranking content in language and from the literature [16, 23]. From the translated corpus, we tagged each token using the Penn Treebank POS[10] tagset built up of 36 different tags. Each of the tags was then mapped into each correspondence group. For example Noun Singular, Noun Plural, Proper Noun Singular and Proper Noun Plural were mapped into the group NOUN using the Universal POS tagset[11]. Only the tokenized words associated with the four selected POS groups, as described above, were retained in the text; all other words were removed. We then used the word-tag pair to represent the syntactic feature as it had been shown to produced better sentiment classification performance when used together [24].

*4) Idiosyncratic Features Extraction:*
Motivated by study in the field of linguistics, we created a syntax rule in the form of NOUN-ADPOSITION-NOUN to identify idiosyncratic phraseology from the corpus. For example, the phrase 'head of cabbage' was identified as an idiosyncratic phrase as the 'head', 'of' and 'cabbage' will be tagged as noun, adposition and noun by the POS tagger. In this work, we are only interested with idiosyncratic phrases in the context of sarcasm detection. As a result, 177 idiosyncratic features were identified from the corpus. Examples of idiosyncratic features found in the corpus were 'puppet of buffalo', 'face of sharks', 'kinds of beans', 'ants than elephants', 'people in clown' and 'joke between continents'. The identified idiosyncratic features were then replaced in the text (for both sarcastic and non-sarcastic examples) with a unique identifier, *idiosyncratic_x*, where $1 \leq x \leq 177$.

## IV. EXPERIMENTS

This section presents the results obtained with respect to the experiments conducted to evaluate the proposed feature extraction process for sarcasm detection.

### A. The Dataset

The bilingual corpus used for evaluation purposes was acquired from comments related to economic news from Facebook public pages. The comments were acquired using Graph API Explorer[12] and the Facebook Query Language (FQL). All links, pictures and video were filtered from the data. Thus only text based comments were considered. A total of 3000 comments were acquired and annotated manually, according to whether they featured sarcasm or not, by three annotators. The annotations were considered to be "valid" only if all three annotators agreed. In this manner a subset of 1970 comments was derived from the original set of 3000 comments, 969 sarcastic comments and 1001 non-sarcastic comments. The Fleiss's kappa inter annotator agreement score for sarcasm was 0.47, which is "moderate agreement" [25].

### B. Experimental Setup

The objective of the evaluation was to identify the effectiveness and the best combination of the different categories of features identified, for the detection of sarcasm. The proposed extracted features (see Section III) were vectorized using Term Frequency - Inverse Document Frequency (TF-IDF) vectorization and normalized to document length. Five sets of experiment were conducted as itemized in Table II. We used the F-measure (Fm) value to evaluate the results. Fm is the harmonic mean of the recorded precision (correct predicted value) and the recall (actual correct prediction). All experiments were conducted using 10-fold cross validation in the Weka Knowledge Flow [26]. The results are presented in the following sections.

TABLE II. EXPERIMENTAL COMBINATION OF FEATURES

| Experiment | |
|---|---|
| Phase I: | Single feature |
| Phase II: | Two combination |
| Phase III: | Three combination |
| Phase IV: | Four combination |
| Phase V: | Five combination (All) |

*1) Feature Selection:*
Pearson's correlation coefficient was used for feature selection. This operates by ranking features according to their correlated class (sarcasm or non-sarcasm). In this paper, we selected the top 25%, 50% and 75% proportion of the features for classifier generation. The aim was to observe the effect of different numbers of features on sarcasm detection performances. We also investigate the performance if without feature selection. The details of the number of features in each category, as used with respect to the experiment reported here, are given in Table III.

*2) Classification and Parameter Setting:*
For the experiments, non-linear SVM was used as the classification model because it had been shown to perform well in the context of similar domains found in the literature [13, 16]. The variation of non-linear SVM used was LibSVM [27] as provided as part of the Weka data mining workbench with default parameters of C = 3.0, Gamma = 0.03 and Radial Basis Function (RBF) as the kernel function.

## V. RESULT AND ANALYSIS

The results of the experiment when using the proposed features to detect sarcasm are shown in Table IV. The best performance for single feature groups was produced using the syntactic category (Fm = 0.847) and 50% of the features. The combinations of syntactic and prosodic recorded the best performance for two combinations with an Fm score of 0.851.

---

TABLE III.    THE NUMBER OF FEATURES USED FOR EXPERIMENTATION

| Experiment | % Feature Selection (FS) size | 25% | 50% | 75% | No FS |
|---|---|---|---|---|---|
|  | Feature | Feature size | | | |
| Phase I | 1. Lexical (Baseline) | 683 | 1365 | 2048 | 2730 |
|  | 2. Pragmatic | 1 | 2 | 3 | 4 |
|  | 3. Prosodic | 17 | 33 | 50 | 66 |
|  | 4. Syntactic | 924 | 1848 | 2771 | 3695 |
|  | 5. Idiosyncratic | 44 | 89 | 133 | 177 |
| Phase II | 6. Lexical + Pragmatic | 684 | 1367 | 2051 | 2734 |
|  | 7. Lexical + Prosodic | 699 | 1398 | 2097 | 2796 |
|  | 8. Lexical + Syntactic | 1606 | 3213 | 4819 | 6425 |
|  | 9. Lexical + Idiosyncratic | 727 | 1454 | 2180 | 2907 |
|  | 10. Syntactic + Pragmatic | 925 | 1850 | 2774 | 3699 |
|  | 11. Syntactic + Prosodic | 940 | 1881 | 2821 | 3761 |
|  | 12. Syntactic + Idiosyncratic | 968 | 1936 | 2904 | 3872 |
|  | 13. Pragmatic + Prosodic | 18 | 35 | 53 | 70 |
|  | 14. Pragmatic + Idiosyncratics | 45 | 91 | 136 | 181 |
|  | 15. Prosodic + Idiosyncratics | 61 | 122 | 182 | 243 |
| Phase III | 16. Lexical + Pragmatic + Prosodic | 700 | 1400 | 2100 | 2800 |
|  | 17. Lexical + Pragmatic + Idiosyncratic | 728 | 1456 | 2183 | 2911 |
|  | 18. Lexical + Prosodic + Idiosyncratic | 743 | 1487 | 2230 | 2973 |
|  | 19. Syntactic + Pragmatic + Prosodic | 941 | 1883 | 2824 | 3765 |
|  | 20. Syntactic + Pragmatic + Idiosyncratic | 969 | 1938 | 2907 | 3876 |
|  | 21. Syntactic + Prosodic + Idiosyncratic | 985 | 1969 | 2954 | 3938 |
|  | 22. Pragmatic + Prosodic + Idiosyncratic | 62 | 124 | 185 | 247 |
| Phase IV | 23. Lexical + Pragmatic + Prosodic + Idiosyncratic | 744 | 1489 | 2233 | 2977 |
|  | 24. Syntactic + Pragmatic + Prosodic + Idiosyncratic | 986 | 1971 | 2957 | 3942 |
| Phase V | 25. Lexical + Syntactic + Pragmatic + Prosodic + Idiosyncratic | 1668 | 3336 | 5004 | 6672 |

The best overall performance was recorded by the combination of the syntactic, pragmatic and prosodic feature categories with an Fm score of 0.852, followed by the four combination of syntactic, pragmatic, prosodic and idiosyncratic with an Fm score of 0.848. However the combination of all features only recorded a best Fm score of 0.825, lower than the baseline score. All the best scores were recorded when retaining 50% of the original features. All the result obtained using feature selection outperformed the results obtained where feature selection was not used.

The idiosyncratic feature did not perform well compared to the other features. It was conjectured that the method used to extract the idiosyncratic features could have affected the performance. It was also conjectured that better result could have been obtained if each phrase had been annotated by

annotators to confirm the peculiar and odd phrase. The creation and use of a Malay idiosyncratic feature category might also improve performance.

The overall result demonstrated the effectiveness of the use of syntactic feature for sarcasm detection compared to lexical features. Combining categories of syntactic features produced better result, in contrast to the use of lexical features when adding more feature categories tended to slightly reduce the performance. The combination of lexical and syntactic categories also produced lower performance, it was conjectured that this was due to the repetition of the same features. Therefore we conclude the best combination for sarcasm detection in the context of Malay social media data was the three combination of syntactic, pragmatic and prosodic feature categories.

TABLE IV.    CLASSIFICATION PERFORMANCES

| Experiment | % Feature Selection (FS) size | 25% | 50% | 75% | No FS |
|---|---|---|---|---|---|
| | Feature | F-measure (Fm) | | | |
| Phase I | 1. Lexical (Baseline) | 0.783 | 0.840 | 0.775 | 0.708 |
| | 2. Pragmatic | 0.397 | 0.407 | 0.407 | 0.407 |
| | 3. Prosodic | 0.556 | 0.566 | 0.567 | 0.560 |
| | 4. Syntactic | 0.761 | **0.847** | 0.724 | 0.656 |
| | 5. Idiosyncratic | 0.388 | 0.427 | 0.461 | 0.461 |
| Phase II | 6. Lexical + Pragmatic | 0.777 | 0.836 | 0.768 | 0.702 |
| | 7. Lexical + Prosodic | 0.782 | 0.835 | 0.764 | 0.705 |
| | 8. Lexical + Syntactic | 0.757 | 0.822 | 0.684 | 0.538 |
| | 9. Lexical + Idiosyncratic | 0.778 | 0.834 | 0.769 | 0.707 |
| | 10. Syntactic + Pragmatic | 0.765 | 0.850 | 0.731 | 0.662 |
| | 11. Syntactic + Prosodic | 0.756 | **0.851** | 0.738 | 0.659 |
| | 12. Syntactic + Idiosyncratic | 0.766 | 0.845 | 0.725 | 0.657 |
| | 13. Pragmatic + Prosodic | 0.586 | 0.594 | 0.598 | 0.596 |
| | 14. Pragmatic + Idiosyncratics | 0.429 | 0.429 | 0.429 | 0.428 |
| | 15. Prosodic + Idiosyncratics | 0.542 | 0.590 | 0.625 | 0.625 |
| Phase III | 16. Lexical + Pragmatic + Prosodic | 0.773 | 0.831 | 0.762 | 0.708 |
| | 17. Lexical + Pragmatic + Idiosyncratic | 0.777 | 0.827 | 0.763 | 0.705 |
| | 18. Lexical + Prosodic + Idiosyncratic | 0.778 | 0.827 | 0.762 | 0.716 |
| | 19. Syntactic + Pragmatic + Prosodic | 0.755 | **0.852** | 0.737 | 0.664 |
| | 20. Syntactic + Pragmatic + Idiosyncratic | 0.767 | 0.844 | 0.733 | 0.659 |
| | 21. Syntactic + Prosodic + Idiosyncratic | 0.770 | 0.846 | 0.732 | 0.662 |
| | 22. Pragmatic + Prosodic + Idiosyncratic | 0.571 | 0.614 | 0.641 | 0.641 |
| Phase IV | 23. Lexical + Pragmatic + Prosodic + Idiosyncratic | 0.770 | 0.827 | 0.761 | 0.707 |
| | 24. Syntactic + Pragmatic + Prosodic + Idiosyncratic | 0.761 | **0.848** | 0.739 | 0.664 |
| Phase V | 25. Lexical + Syntactic + Pragmatic + Prosodic + Idiosyncratic | 0.765 | **0.825** | 0.674 | 0.533 |

## VI. CONCLUSION

An approach to detecting sarcasm in bilingual texts using combinations of categories of features extraction using NLP has been proposed. The approach extracts the features from corpuses in either bilingual or translated form. Five categories of NLP feature were considered: lexical, pragmatic, prosodic, syntactic and idiosyncratic. A non-linear SVM was used for classification purposes with respect to sarcasm detection, to evaluate the feature categories (used on their own and in combination). Comparison with a baseline feature demonstrated that the proposed approach performed better.

For future work, the intention is to build a framework for sarcasm detection and classification using the proposed features to support sentiment analysis.

## ACKNOWLEDGMENT

## REFERENCES

[1]    A. Farzindar and D. Inkpen, *Natural Language Processing for Social Media* vol. 8: Morgan & Claypool Publishers, 2015.

[2]    B. Liu, *Sentiment analysis: Mining opinions, sentiments, and emotions*: Cambridge University Press, 2015.

[3]    N. Indurkhya and F. J. Damerau, *Handbook of natural language processing*, Second Edition ed.: CRC Press, 2010.

[4]    A. Reyes, P. Rosso, and D. Buscaldi, "From humor recognition to irony detection: The figurative language of social media," *Data & Knowledge Engineering,* vol. 74, 2012, pp. 1-12.

[5] F. Provost and T. Fawcett, *Data Science for Business: What you need to know about data mining and data-analytic thinking*: " O'Reilly Media, Inc.", 2013.

[6] R. Gibbs and H. Colston, "The future of irony studies," in *Irony in language and thought: A cognitive science reader*, R. Gibbs and H. Colston, Eds., ed London: Taylor & Francis Group, 2007, pp. 339-360.

[7] R. J. Kreuz and R. M. Roberts, "On satire and parody: The importance of being ironic," *Metaphor and Symbolic Activity,* vol. 8, 1993/06/01 1993, pp. 97-109.

[8] R. W. Gibbs, "Irony in talk among friends," *Metaphor and symbol,* vol. 15, 2000, pp. 5-27.

[9] M. L. Dress, R. J. Kreuz, K. E. Link, and G. M. Caucci, "Regional variation in the use of sarcasm," *Journal of Language and Social Psychology,* vol. 27, 2008, pp. 71-85.

[10] J. Daniel and H. James, "Speech and Language processing: An introduction to natural language processing," *Computational Linguistics and Speech Recognition, 2nd Ed., Prentice Hall,* 2009.

[11] D. Bikel and I. Zitouni, *Multilingual Natural Language Processing Applications: From Theory to Practice*: IBM Press, 2012.

[12] G. Lakoff and M. Johnsen, "Metaphors we live by. London: The university of Chicago press," *Prieiga per internetą:* http://shu. bg/tadmin/upload/storage/161. pdf [žiūrėta 2012 09 24], 2003.

[13] E. Lunando and A. Purwarianti, "Indonesian social media sentiment analysis with sarcasm detection," in *2013 International Conference on Advanced Computer Science and Information Systems (ICACSIS)*, 2013, pp. 195-198.

[14] C. C. Liebrecht, F. A. Kunneman, and A. P. J. van den Bosch, "The perfect solution for detecting sarcasm in tweets# not," in *4th Workshop on Computational Approaches to Subjectivity, Sentiment and Social Media Analysis*, Atlanta, Georgia, 2013, pp. 29-37.

[15] N. Littlestone, "Learning quickly when irrelevant attributes abound: A new linear-threshold algorithm," *Machine learning,* vol. 2, 1988, pp. 285-318.

[16] T. Ptácek, I. Habernal, and J. Hong, "Sarcasm detection on Czech and English Twitter," in *COLING 2014 - 25th International Conference on Computational Linguistics*, Dublin, Ireland, 2014, pp. 213-223.

[17] B. Charalampakis, D. Spathis, E. Kouslis, and K. Kermanidis, "A comparison between semi-supervised and supervised text mining techniques on detecting irony in Greek political tweets," *Engineering Applications of Artificial Intelligence,* 2016.

[18] S. Agarwal, S. Godbole, D. Punjani, and S. Roy, "How much noise is too much: A study in automatic text classification," in *Data Mining, 2007. ICDM 2007. Seventh IEEE International Conference on*, 2007, pp. 3-12.

[19] G. Forman, "Chapter: Feature Selection for Text Classification Book: Computational Methods of Feature Selection Chapman and Hall/CRC Press, 2007," 2007.

[20] P. Carvalho, L. Sarmento, M. J. Silva, and E. de Oliveira, "Clues for detecting irony in user-generated contents: Oh...!! it's "so easy" ;-)," in *Proceeding of the 1st international CIKM workshop on Topic-sentiment analysis for mass opinion - TSA '09*, New York, New York, USA, 2009, p. 53.

[21] A. Balahur and M. Turchi, "Comparative experiments using supervised learning and machine translation for multilingual sentiment analysis," *Computer Speech & Language,* vol. 28, 1// 2014, pp. 56-75.

[22] D. Castelvecchi. (2016, 23/1/2017). Deep learning boosts Google Translate tool. Available: http://www.nature.com/news/deep-learning-boosts-google-translate-tool-1.20696

[23] O. Jespersen, *The philosophy of grammar*: University of Chicago Press, 1992.

[24] R. Xia and C. Zong, "Exploring the use of word relation features for sentiment classification," in *Proceedings of the 23rd International Conference on Computational Linguistics: Posters*, 2010, pp. 1336-1344.

[25] J. R. Landis and G. G. Koch, "The measurement of observer agreement for categorical data," *biometrics,* 1977, pp. 159-174.

[26] M. Hall, E. Frank, G. Holmes, B. Pfahringer, P. Reutemann, and I. H. Witten, "The WEKA data mining software: an update," *ACM SIGKDD explorations newsletter,* vol. 11, 2009, pp. 10-18.

[27] C.-C. Chang and C.-J. Lin, "LIBSVM: A library for support vector machines," *ACM Transactions on Intelligent Systems and Technology (TIST),* vol. 2, 2011, p. 27.