

DWAEF: a deep weighted average ensemble framework harnessing novel indicators for sarcasm detection

Richa Sharma^{a,*}, Simrat Deol^b, Udit Kaushish^c, Prakher Pandey^d and Vishal Maurya^e

^a *Department of Computer Science, Keshav Mahavidyalaya, University of Delhi, India*
E-mail: richasharma@keshav.du.ac.in; ORCID: <https://orcid.org/0000-0002-4472-1681>

^b *Department of Computer Science, Keshav Mahavidyalaya, University of Delhi, India*
E-mail: simrat205711@keshav.du.ac.in

^c *Department of Computer Science, Keshav Mahavidyalaya, University of Delhi, India*
E-mail: udit205805@keshav.du.ac.in

^d *Department of Computer Science, Keshav Mahavidyalaya, University of Delhi, India*
E-mail: prakher205723@keshav.du.ac.in

^e *Department of Computer Science, Keshav Mahavidyalaya, University of Delhi, India*
E-mail: vishal205750@keshav.du.ac.in

Abstract. Sarcasm is a linguistic phenomenon often indicating a disparity between literal and inferred meanings. Due to its complexity, it is typically difficult to discern it within an online text message. Consequently, in recent years sarcasm detection has received considerable attention from both academia and industry. Nevertheless, the majority of current approaches simply model low-level indicators of sarcasm in various machine learning algorithms. This paper aims to present sarcasm in a new light by utilizing novel indicators in a deep weighted average ensemble-based framework (DWAEF). The novel indicators pertain to exploiting the presence of simile and metaphor in text and detecting the subtle shift in tone at a sentence's structural level. A graph neural network (GNN) structure is implemented to detect the presence of simile, bidirectional encoder representations from transformers (BERT) embeddings are exploited to detect metaphorical instances and fuzzy logic is employed to account for the shift of tone. To account for the existence of sarcasm, the DWAEF integrates the inputs from the novel indicators. The performance of the framework is evaluated on a self-curated dataset of online text messages. A comparative report between the results acquired using primitive features and those obtained using a combination of primitive features and proposed indicators is provided. The highest accuracy of 92% was achieved after applying DWAEF, the proposed framework which combines the primitive features and novel indicators together as compared to 78.58% obtained using Support Vector Machine (SVM) which was the lowest among all classifiers.

Keywords: sarcasm detection, deep ensemble learning, weighted average ensemble model, graph neural networks, BERT, fuzzy logic

1. Introduction

Natural languages have evolved gracefully over time all around the globe. Various nuances of a language allow humans to put forth their views on myriad topics with ease and creativity. The use of

*Corresponding author. E-mail: richasharma@keshav.du.ac.in.

figurative language by native speakers is one such medium of expressing opinions [1]. Sarcasm, interlaced with irony and wit, affords both sharpness and subtlety to convey contempt. Automatic detection of sarcasm in text is one of the critical challenges faced by researchers in the field of sentiment analysis. Sensing the negative connotation in a sentence containing positive words is required to detect sarcasm in an effective manner. Moreover, sarcasm is closely tied to linguistic and cultural norms and practices, and its use can vary significantly between different languages and cultures. While sarcasm has clear definitions in the literature, its interpretation may be greatly affected by the cultural background and contextual knowledge of the reader [2] [3]. Sarcasm can be considered rude or disrespectful in some cultures while being accepted as a form of humour in others. Primitive computational models developed for sarcasm detection made use of primitive features such as n-grams, punctuation and intensifiers and exploited machine learning algorithms for classification purposes.

To identify sarcasm in text in an automated manner, the present study proposes a deep weighted average ensemble-based framework (DWAEF). The proposed framework makes use of three indicators to produce competent results. These indications concern utilising the presence of simile and metaphor in text and identifying small shifts in tone between the constituent clauses of a sentence. The framework leverages deep learning components, namely graph neural network (GNN) [4] and bidirectional encoder representations from transformers (BERT) [5] based embeddings to detect simile and metaphor respectively and fuzzy logic [6] to apprehend polarity shifts between the constituent clauses of a sentence. Finally, the outputs of the three components are provided to DWAEF, an ensemble structure comprising attentive interpretable tabular learning (TabNet) [7], one-dimensional convolutional neural networks (1-D CNN) and multilayer perceptron (MLP) based learners. The results obtained using the ensemble method are thoroughly compared with results obtained using only the base learners. With the accuracy of 92.01% achieved by DWAEF, the proposed ensemble-based approach surpasses the results obtained based on the usage of primitive features only. The main contributions of this study are summarized below:

- (1) Leveraging key linguistic features, namely- simile, metaphor and constituent clauses of a sentence for sarcasm detection that, to the best of the authors' knowledge, have not yet been used together for this purpose
- (2) Implementing GNN in the framework to detect the presence of simile in a text on the basis of a sentence's dependency tree
- (3) Exploiting BERT embeddings to detect the presence of metaphor in a text
- (4) Capturing the shift in polarity of a sentence's constituent clauses using fuzzy-logic
- (5) Harnessing ensemble structure of various deep learning algorithms for facilitating sarcasm classification tasks.

The rest of the paper is organized as follows. Section 2 discussed the earlier research done in the field of sarcasm detection. Section 3 puts forth the motivation behind the proposed study. Section 4 presents and describes the proposed methodology. Section 5 describes the experiments and gives a detailed analysis of the results obtained. Section 6 concludes the paper.

2. Related work

Detection of sarcasm is a challenge for humans and for machines, even more so. As a result, it has gained popularity in many natural language processing (NLP) applications such as social media analysis, sentiment analysis and customer service [8]. The study in [9] contended that sarcasm is distinct from

1 other forms of figurative language and contributed to the understanding of the linguistic and cognitive 1
2 processes underlying sarcasm. An extensive survey of the literature brought to light that researchers have 2
3 mainly employed machine learning techniques to a variety of features, including lexical, syntactic, and 3
4 semantic features for detecting sarcasm. The common form of sarcasm consists of a positive sentiment 4
5 situation followed by a negative sentiment situation. The study in [10] discussed an algorithm that auto- 5
6 matically learns positive and negative sentiment phrases from sarcastic tweets. While the approach was 6
7 innovative and showed promising results, there were several potential drawbacks to this method related 7
8 to limited applicability, lack of context sensitivity and limited evaluations. In [11] authors surveyed sev- 8
9 eral machine learning algorithms to classify the sarcastic tweets and found that a combination of Support 9
10 Vector Machine (SVM) and convolutional neural network (CNN) resulted in higher prediction accuracy. 10
11 The review also covered various aspects of sarcasm detection, including datasets used, machine learning 11
12 algorithms employed, feature selection and evaluation metrics. The authors found that while there was 12
13 a growing interest in using machine learning algorithms for sarcasm detection in Twitter, the field still 13
14 faced several challenges, such as the lack of annotated data and the difficulty of detecting sarcasm in 14
15 short, informal messages. Researchers in [12] applied K-Nearest Neighbours (KNN), Random Forest 15
16 (RF), SVM and max entropy techniques on the following features- sentiment related, syntactic and se- 16
17 mantic, punctuation-related and pattern-related. As per their results, the RF classifier outperformed all 17
18 applied models with the highest accuracy of 83.1%. Their approach relied heavily on pre-defined pat- 18
19 terns, which failed to cover all instances of sarcasm on Twitter. New patterns might need to be added or 19
20 existing patterns might need to be modified as language and culture evolve over time. 20
21

22 In [13], the researchers harvested sarcastic tweets with the help of hashtags such as *#not*, *#sarcasm* 22
23 and divided them into tweets containing user mentions and tweets that do not. A trained machine learn- 23
24 ing classifier, Winnow2 [14] was then employed to segregate tweets aimed at specific users from those 24
25 that were not. This study focused on a specific type of social media messaging, which might not have 25
26 been representative of all types of social media discourse. The linguistic markers identified in this study 26
27 might not have been applicable to other types of messages or platforms. Additionally, the research relied 27
28 on a limited dataset of social media messages, potentially introducing bias in the analysis and limiting 28
29 the generalizability of the results. In [15], the researchers included extra-linguistic information and em- 29
30 ployed binary logistic regression with l2 regularization and achieved a gain in accuracy as compared 30
31 to purely linguistic features in sarcasm detection. The paper relied heavily on contextual information. 31
32 While the approach had some strengths, such as the use of advanced machine learning techniques, it also 32
33 had potential drawbacks, including limited evaluation, over-reliance on context, complexity and limited 33
34 generalizability. In [16], authors used unigrams, bigrams and trigrams to create more general sarcasm 34
35 indicators which in turn resulted in a 75% precision and 62% recall score for a bootstrapping classifier. 35
36 However, since their bootstrapping approach relied on the performance of existing classifiers, it might 36
37 not prove to be optimal for sarcasm and nastiness detection. If the initial classifier is not accurate, the 37
38 bootstrapping approach might not improve its performance significantly. The authors of [17] presented 38
39 a case study on the identification of nonliteral language, specifically sarcasm, in social media. The au- 39
40 thors analyzed different approaches to detecting sarcasm in Twitter data, including sentiment analysis, 40
41 rule-based techniques and machine learning. They found that a combination of these methods was the 41
42 most effective in identifying sarcasm and that the accuracy of the detection was influenced by the context 42
43 and topic of the conversation. This study's limitations include its Twitter-specific focus, limited scope 43
44 regarding context and lack of comprehensive analysis regarding different factors that may affect sarcasm 44
45 detection, such as user demographics. 45
46

In recent years there has also been a shift towards using neural networks due to their superior performance in a variety of machine learning tasks. Other machine learning techniques, such as Decision Trees or SVM, have limitations in their ability to handle complex and high-dimensional data, such as language data and may require extensive feature engineering. Neural networks excel at recognizing and identifying patterns in large and complex datasets, including language data. By training on a sizeable dataset of sarcastic and non-sarcastic language, neural networks can detect the distinct features of sarcasm, such as negation, exaggeration and tone of voice. Therefore, they are an ideal tool for sarcasm detection, as they can learn from vast amounts of data and identify subtle patterns that other machine learning algorithms may struggle to recognize. Researchers in recent studies have employed various neural network techniques such as CNN and long short-term memory (LSTM) along with different word embeddings namely, Word2Vec, FastText and GloVe on the Reddit Corpus. They accounted for the impact of varying epochs, training size and dropout on the performance [18–20]. These studies only used word embeddings and hyperparameter tuning as features to identify sarcasm, which might not be sufficient to capture all the nuances of sarcasm in language. In [21], the study implemented BERT, robustly optimized BERT pretraining approach (ROBERTa), LSTM, bi-directional long short-term memory (Bi-LSTM) and bi-directional gated recurrent unit (Bi-GRU) models for detecting sarcasm in text. They concluded that the transformer-based ensemble performed better than the baseline models scoring 0.43 on F1-score. Authors in [22] used an ensemble of LSTM, gated recurrent unit (GRU) and Baseline CNNs to detect sarcasm in online text and concluded using a weighted average ensemble resulted in better results. However, the approach used by the researchers failed to detect sarcastic tweets written in a very polite way. In [23], the researchers used four component methods namely LSTM, CNN-LSTM, SVM and MLP on the Reddit and Twitter datasets resulting in F1 scores of 67% and 73% respectively.

Recent studies have made extensive use of word embeddings in deep neural networks for various NLP tasks. However, there is a growing demand for modelling text data as graphs. In comparison to revolutionary neural networks such as recurrent neural networks (RNN), LSTM, CNN and BERT, the graphical representation of text allows for more efficient extraction of semantic and structural information. Therefore, numerous researchers have investigated graph-based methods and their application to NLP problems. The first graph attention-based model to identify sarcasm on social media was proposed in [24]. The graph model captured complex relationships between a sarcastic tweet and its conversational context by modelling a user’s social and historical context together. GNN, a modern class of networks applied on graph-structured data [25], has found application in the field of sarcasm detection. In [26], a graph convolutional network (GCN) was used to capture global information features in a satirical context and a Bi-LSTM was implemented to capture the sequence features of the comments. The two sets of results were combined and evaluated using a conventional classifier. Later, The authors in [27] proposed an affective dependency graph convolutional network framework to detect messages with implied contradictions and incongruity. Apart from the state-of-the-art technology, many researchers have investigated the use of different features in text data and different methodologies for detecting sarcasm. The article [28] provides a detailed literature survey on sarcasm detection. Additionally, it provided a detailed analysis of the set of features used for sarcasm detection.

Subsection 2.1 discusses the types of feature sets used in sarcasm detection and how researchers have employed them in past studies.

2.1. Types of primitive feature sets used in sarcasm detection

Previous works on sarcasm detection made use of low-level features such as n-grams, punctuation, intensifiers and so forth. Some of the primitive features such as punctuation count, count of mixed-case

words, count of repeated words and letters, presence of intensifiers and presence of interjections are later used in this research as part of feature set preparation. The primitive features used in sarcasm detection can be broadly classified into:

- (1) Lexical features: This feature set includes text properties such as unigrams, bigrams, n-grams, skip-gram, hashtags, etc. The study in [16] used unigrams, bigrams and trigrams to create more general sarcasm indicators thereby improving the precision and recall of their bootstrapping classifier. Researchers in [15] created binary indicators of lower-cased word unigrams and bigrams along with brown cluster unigrams and bigrams which grouped words used in similar contexts into the same cluster. The authors of [10] extracted every unigram, bigram and trigram that occurred immediately right after a positive sentiment phrase in a sarcastic tweet.
- (2) Paralinguistics features: These are some of the main features used for sarcasm detection in text. They include emoticons, smileys, number of hashtags, replies and so forth. The study in [12] included the count of positive, negative and sarcastic emoticons. The authors of [29] considered the effect of sentiment contained in hashtags by developing a set of rules around the number of hashtags and their polarity. Researchers in [17] took into account the sentiment of replies to the user.
- (3) Hyperbole features: These features include intensifiers, interjections, quotes, punctuation and so forth. Researchers in [15] created a binary indicator for the presence of 50 intensifiers retrieved from Wikipedia. The authors of the study in [30] opined that writers often use sarcasm-based writing styles to compensate for the lack of visual or verbal cues. The authors in [31] accounted for uppercase and lowercase characters along with the repetition of punctuation marks.
- (4) Contextual features: These features comprise extra components, outside the realm of formal linguistics, used frequently in a sentence, especially in online messages. The researchers in [13] harvested a large number of sarcastic tweets with the help of hashtags such as #not#sarcasm and divided them into tweets containing user mention and tweets that do not. In [15], the researchers emphasized extra-linguistic information from the context of a tweet in the form of ‘author features’, ‘audience features’, ‘environment features’ and ‘tweet features’ and achieved gains in accuracy compared to purely linguistic features in sarcasm detection.

The previous research and development in the field of sarcasm detection prompted the authors of this paper to take on this problem and address its concerns at a new level. The following section describes the impetus behind the present study.

3. Motivation

This section explains the rationale for delving into the complexities of sarcasm detection using similes, metaphors and the clausal structure of a sentence. Subsection 3.1 discusses similes in literature and forms the base for the proposed methodology for its computational detection. Subsection 3.2 introduces and explores metaphors in literature thereby building the foundation for its computational detection. Subsection 3.3 deliberates upon a sentence’s clausal structure as well as the polarity change from one clause to another. 3.4 lays the motivation for using deep learning methods in an ensemble structure.

3.1. Simile

A simile is a figurative device used to draw comparisons between two unlike things. Its presence is always explicitly indicated with the usage of “like” or “as”. A simile consists of the following four key

Table 1
Examples and constituent components of a simile

Simile	Tenor	Vehicle	Property	Event
Her voice is as smooth as silk.	Her voice	silk	Smoothness	is
A sweet voice carolling like a gold-caged nightingale	sweet voice	gold-caged nightingale	The property here is implicit, left for the reader to infer	carolling
Her grandmother's love story was as old as the hills.	grandmother's story	love hills	old	was
A slow thought that crept like a cold worm through his brain.	slow thought	cold worm	The property here is implicit, left for the reader to infer	crept

components- *tenor*, *vehicle*, *property*, *event* and *comparator* [32]. Table 1 provides examples of similes along with their constituent components. This study proposes the presence of a simile as a potential marker for sarcasm in the text as its presence in a sarcastic remark may accentuate the hidden emotion. For instance, “*Of course they were invited! They are always as welcome as a skunk at a lawn party*” implies that the subject’s presence is actually not appreciated. Here the comparison “*as welcome as a skunk at a lawn party*” represents the undesirability and vileness of the subject. Another potential example of a sarcastic remark embedding a simile is “*Asking politicians to give up a source of money is like asking Dracula to forsake blood*” wherein the speaker mocks politicians’ flaws by drawing analogies to Dracula. The computational detection of simile relies on the syntactic dependency tree of a sentence, which is described in more detail in section 4.

3.2. Metaphor

A metaphor is a figure of speech that compares two unrelated ideas. At the basic linguistic level, both metaphor and simile involve the juxtaposition of two concepts. However, metaphors lack the usage of “*like*” or “*as*” while drawing the comparison. For example, the two statements, “*Mary is a rock*” and “*Mary is like a rock*” will be inferred by the reader in the same sense about Mary’s personality [33]. The only difference between the statements is that the former statement is a metaphor and the latter is a simile. The difference lies in the presence of the comparator “*like*” in one and its absence in the other.

A metaphor also arises when seemingly unrelated properties of one concept are seen in terms of the properties of some other concept. Metaphorical utterances in sarcastic remarks in certain situations are common. For example, “*You are the cream in my coffee*” when used sarcastically implies that the hearer is an unpleasant addition to the speaker’s life [34]. Another example of such an utterance is, “*I am not saying that I hate you, what I am saying is that you are literally the Monday of my life.*” wherein the speaker indirectly expresses his hate towards the listener by comparing the latter’s presence in the former’s life as depressing and unwanted as Monday. Since comparison is drawn between two distinctive entities, computing cosine similarity between the subject and object of comparison forms the bases for its computational detection. This study facilitates the detection of only two types of metaphorical sentences out of the three mentioned by [35]. Table 2 provides a summary of the two types. The third type of metaphor takes the form of an Adjective acting on a Noun, for example, “*He has a fertile imagination*”. Due to the different linguistic structure of adjectival metaphors a separate analysis and methodology is required for its detection, which is beyond the scope of this current paper.

Table 2
Types of metaphors addressed in this study and their examples

Metaphor type	Relationship	Example
Type I	Subject IS A object phrase; (X is Y)	1. Mary is a rock. 2. He is the sugar in my coffee. 3. That fella is the raspberry seed in my wisdom tooth.
Type II	Verb acting on Noun phrase; (X acts on Y)	1. My car drinks gasoline. 2. He planted good ideas in their minds. 3. Inflation has eaten up all my savings.

Table 3
Distinctive subtleties between main and subordinate clauses

Sentence	Main clause	Subordinating conjunction	Subordinate clause
She had a long career but she is remembered for one early work.	She had a long career	but	-
I first saw her in Paris, where I lived in the early nineties.	I first saw her in Paris	where	(where) I lived in the early nineties
If it looks like rain, a simple shelter can be made out of a plastic sheet.	A simple shelter can be made out of the plastic sheet	if	(if) it looks like rain

3.3. Clauses

Clauses are a group of related words which unlike phrases have a subject and a verb. A clause can be a part of a sentence or be a complete sentence in itself. All sentences have at least one main clause. The main clause is a clause that can stand alone as an independent complete sentence. On the other hand, a subordinate clause is a clause that cannot stand as an independent complete sentence by itself. It is typically introduced with a subordinating conjunction and is dependent on the main clause. Consider the examples taken from an article¹ given in table 3 elaborating the distinctive subtleties between a main clause and a subordinate clause. Sarcasm, in its prevalent form, exists as the disparity of sentiments. This disparity can further take up two forms [14]:

- (1) A shift from positive polarity to negative polarity: In this type, sarcastic sentences contain positive expressions followed by negative expressions. Consider the sarcastic sentence, “*Thank you, officer; now that you have my license I can’t drive*” where the main clause “*Thank you officer*” has a positive connotation and the subordinate clause “*now that you have my driving license I can’t drive*” has a negative connotation.
- (2) A shift from negative polarity to positivity polarity: In this type, sarcastic sentences contain negative expressions followed by positive expressions. For instance, “*I hate my sister because she cooks so well*” wherein the main clause “*I hate my sister*” holds a negative connotation and the subordinate clause “*because she cooks so well*” holds a positive connotation.

To cater to such types of situations, this research proposes to measure the polarity shift from the main clause of a sentence to the subordinate clause of a sentence at various degrees as a potential indicator of

¹<https://www.lexico.com/grammar/clauses>

sarcasm. The significance of fuzziness comes into play while dealing with ambiguities surrounding how positive or negative a stand-alone clause can be. Its amalgamation with the computational detection of polarity shift may result in efficient results.

3.4. Motivation behind using a deep ensemble structure

Current cutting-edge research studies utilise geometric deep learning, BERT and fuzzy logic. This study combines these techniques into a single framework in order to produce competent results. Furthermore, most authors have employed conventional machine learning classifiers for evaluation purposes, whereas ensembles of deep learning algorithms (TabNet, CNN and MLP) are employed in this research. One of the most significant issues with conventional machine learning techniques is that they frequently fail to capture the underlying characteristics and structure of the data. Consequently, poor performance is observed when these algorithms are applied to datasets that are highly imbalanced, high-dimensional and noisy. [36]. Therefore, it is essential to construct an efficient model, particularly for complex tasks such as sarcasm detection. Ensemble learning is one of the approaches. Ensemble learning strategies enhance the performance and accuracy of a predictive model by merging multiple models. This approach involves training several models, each with unique algorithms or hyperparameters and then fusing their predictions in a manner that maximizes the final outcome. Any ensemble framework comprises a collection of base learners and meta-learners. Base learners, also known as weak learners, are machine learning classifiers whose predictions are combined with those of other weak learners to compensate for their weaknesses. The meta learner or strong learner is the combined learnt model. The promising results obtained by past researchers with different ensemble structures for sarcasm detection motivated the authors of this work to implement a deep ensemble framework DWAEF. The framework is comprehensively described in the forthcoming section.

4. Methodology

The methodology followed by this research is elaborated in Fig. 1. A detailed description of dataset preparation is provided in subsection 4.1. Once the data has been collected, it goes through several stages of preprocessing as elaborated in subsection 4.2. The pre-processed data is then annotated by four expert linguists with 100% agreement that the dataset consisted of both sarcastic and non-sarcastic sentences. Following this, the pre-processed and correctly labelled data goes through the feature extraction process wherein along with the indicators proposed by this research, primitive features are also extracted. Followed by feature set preparation, the results are obtained using the ensemble framework. Each module portrayed in Fig. 1 is explicated in the forthcoming subsections. Subsection 4.3 discusses the primitive feature set preparation. Subsection 4.4 discusses the detection of novel features namely, simile, metaphor and change in the polarity of a sentence's constituent clauses.

4.1. Data set preparation

Description of the dataset: The authors of this work prepared a dataset of 2,891 sentences written in English. The dataset utilized in this study was sourced from various platforms, including Twitter, News Headlines and the SARC datasets. TextBlob [37] was used to check whether the text was written only in English. Out of these, 1,538 were sarcastic and were compiled from various sources- i) 520 sentences were extracted from Twitter with hashtags- #sarcasm, #not, #sarcastic, #irony, #satire between the time

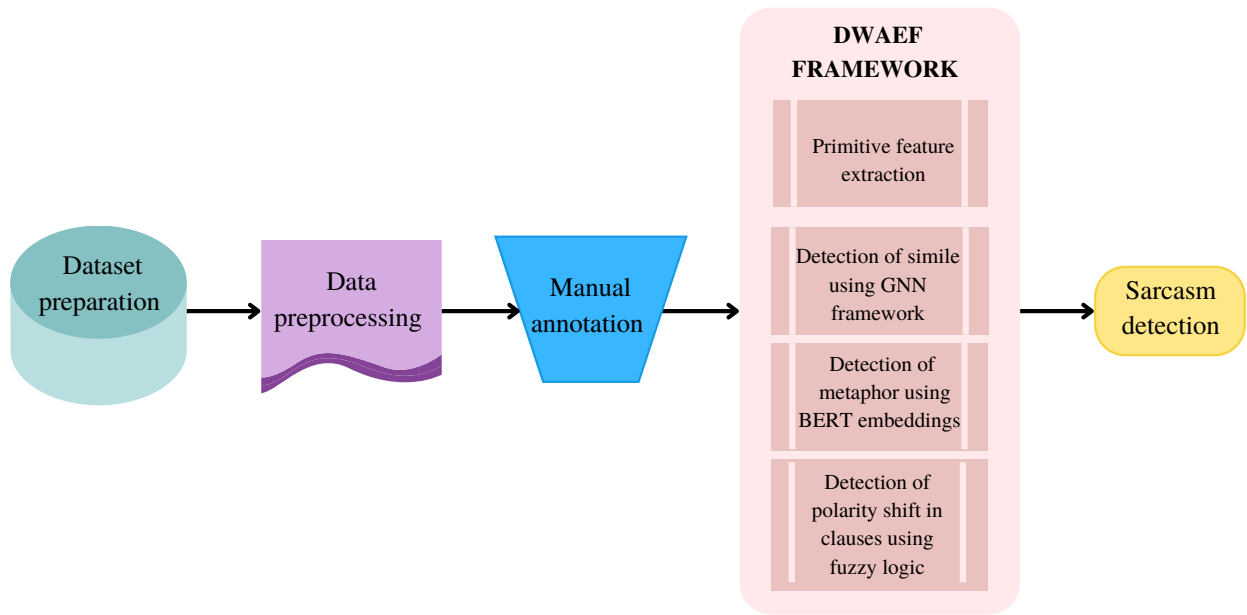


Fig. 1. The methodology of the proposed study.

Table 4
Description of the dataset

Label	SARC	News Headlines	Twitter
sarcastic	498	520	520
non-sarcastic	-	676	677

period of June 2022-October 2022; ii) 520 were taken from the News Headlines dataset curated by [38]; iii) remaining 498 were taken from the SARC dataset curated by [39]. The 1,353 non-sarcastic sentences were compiled from Twitter and the News Headlines dataset. These sources do not belong to the same domain or topic. Twitter data, for instance, encompasses a broad range of subjects, whereas News Headlines may concentrate on specific areas, such as politics, sports, or entertainment. The domain can significantly affect the outcomes of sarcasm detection because the use of sarcasm can vary depending on the context and subject matter. For example, the prevalence of sarcasm in political discussions may differ from its occurrence in conversations about sports or entertainment. The motivation for combining data from various online sources such as Twitter, Reddit and News Headlines is to expose the model to different writing styles as well as myriad domains prevalent online. Moreover, the proposed features in this study operate independently of the domain. Hence, there are no visible effects on the results obtained. This approach is therefore expected to help the model make accurate predictions for a wide range of text messages.

Annotation: Double annotation was performed by four expert linguists who independently performed annotation on the dataset. Once all four annotators had completed their annotations, the annotations were compared with each other to identify any discrepancies or disagreements. A consensus-based approach was followed to resolve these discrepancies or disagreements. This involved having the annotators dis-

cuss and come to a consensus on the correct annotation together. The agreed-upon annotations from all four annotators were used to create a final, consensus annotation for the entire dataset. Further, a series of preprocessing measures were taken. After preprocessing the dataset was reduced to 2,889 sentences. The composition of the dataset and the distribution of sarcastic and non-sarcastic sentences are illustrated in table 4.

4.2. Data preprocessing

The use of slang, hashtags, emoticons, alterations in spelling, and loose punctuation are not inherently redundant and serve various purposes, such as expressing emotions or emphasizing a point. However, they can pose challenges for tasks related to NLP and text analysis, as they often deviate from standard grammar and syntax. The study in [40] provides a more nuanced understanding of how these features of informal text can be modelled and analyzed to improve the quality of NLP frameworks and crowd-sourced annotation tasks. As a result, the authors have accordingly taken these features into account before completely removing them. Certain data pre-processing measures were carried out as follows:

- (1) Duplicate tweets and re-tweets were also dropped.
- (2) Hashtags were completely removed.
- (3) Tweets containing URLs were dropped.
- (4) Emojis were counted and then removed from the text.
- (5) The occurrences of the following punctuation marks (‘.’, ‘?’, ‘*’, ‘!’, ‘,’) were first counted and then the data was freed of irrelevant punctuation marks.

4.3. Primitive features

The primitive features used by this study include various features explained earlier in section 2. The said feature set consists of punctuation count, count of mixed-case words, count of repeated words and letters, presence of intensifiers and presence of interjections. Each one of the preceding features is comprehensively explained below.

- (1) Punctuation Count: The punctuation marks are sometimes overused to indicate sarcasm. For example, to emphasise a point, users use an asterisk (*). To represent a pause, an ellipsis (‘...’) is used and a bunch of exclamatory marks (‘!!!’) indicate exclamatory utterances [30]. Thus, each of the previously described punctuation marks along with some more (‘.’, ‘?’, ‘*’, ‘!’, ‘,’) were counted as one of the features.
- (2) Count of mixed-case words: This feature set includes counting the occurrence of mixed-case words in the text.
- (3) Count of repeated words and letters: Users also tend to repeat letters in words to over-emphasize parts of the text. A similar pattern can be observed in the case of words. As a result, the number of repeated letters and repeated words were counted and used as a set of 2 individual features.
- (4) Presence of intensifiers: Intensifiers or hyperbolic words are generally adverbs or adjectives which strengthen the evaluative utterance of a sarcastic remark. Consider the utterances were taken from [41], “‘fantastic weather’, ‘when it rains’” and “‘weather is good when it rains’”. Both utterances may literally convey a positive outlook of the speaker. However, sensing the context, the utterance with the word fantastic can easily be identified as sarcastic. For this study, a list of commonly used

intensifiers given in appendix A.1 was retrieved from Wikipedia² and used to check the presence of intensifiers in the tweets.

- (5) Presence of interjections: Interjections are words or phrases primarily used in a sentence to convey emotions. For instance, “aha”, “yay”, “oh”, “nah”, “yeah”, “wow” and so forth are some of the commonly used interjections. A list of interjections given in appendix A.2 was retrieved from an article³ and was used to check the presence of interjections in tweets.
- (6) The number of times words having opposite polarities come together: This feature captures the contrast between two words having opposite polarities.
- (7) Length of the largest sequence of words with polarities unchanged
- (8) Count of positive and negative words

4.4. Frameworks for the proposed features

4.4.1. GNN framework for simile detection

For the purpose of this research, simile is detected on the basis of its syntactic pattern using a GNN. The process typically involves first parsing the sentence using a syntactic parser to identify the syntactic relationships between the words. In this case, this was done using Stanford NLP Group’s CoreNLP server [42]. Fig. 2 presents dependency trees of two sentences containing similes created using Stanford NLP Group’s CoreNLP server. The resulting parse tree is then converted into a graph representation. To capture more complex syntactic relationships between the words in a sentence, the authors of this research implemented Bi-Fuse GraphSAGE [43]. To create a representation of each node, GraphSAGE aggregates information from its neighbouring nodes in the graph. This representation is then refined by Bi-Fuse GraphSAGE, which encodes the graph structure in a bi-directional manner by passing messages between nodes in both forward and backward directions, enabling more complex relationships between nodes to be captured. Subgraphs corresponding to specific syntactic structures of a simile are identified by analyzing these refined node representations. Additionally, a graph-level representation is created by combining these refined node embeddings, summarizing the properties of the entire graph. The class label of a sentence is predicted using these graph-level representations. A GNN-based text classification model is used to learn the dependency structure of similes. The entire set-up for the simile classification model consists of a graph construction module, a graph embedding module and a prediction module. Each of the modules is implemented using Graph4NLP library [44]. The modules are elaborated thoroughly below and the entire framework is summarized in Fig. 3.

- (1) **The graph construction module:** The graph construction module focuses on building a syntactic dependency tree-based static graph for each of the texts in the dataset. All of the dependency trees are built using Stanford NLP Group’s CoreNLP server. Dependency relations from the dependency parsing trees are converted into dependency graphs. 3000 sentences consisting of both similes and non-similes were collected for pretraining the GNN framework. Pretraining is done to improve the framework’s ability to recognize similes in the combined dataset. Fig. 4 illustrates a series of initialization preprocessing steps that the raw data goes through before being passed to the graph embedding module. The dataset description, link to the dataset and the initialization preprocessing steps can be found in appendix B.1. Once the initialization is complete, the data items are collated into the batch data which is then used for runtime iteration over the entire dataset.

²<https://en.wikipedia.org/wiki/Intensifier>

³<https://www.english-grammar-revolution.com/list-of-interjections.html>

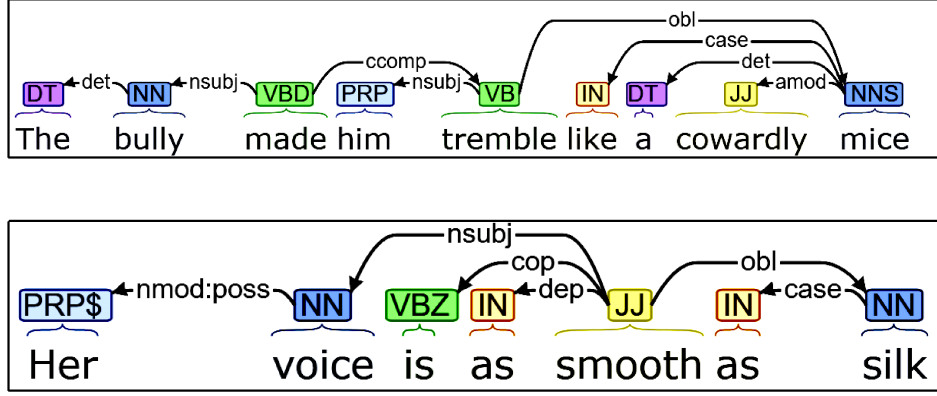


Fig. 2. The dependency trees depicting the syntactic dependency between various components of a simile.

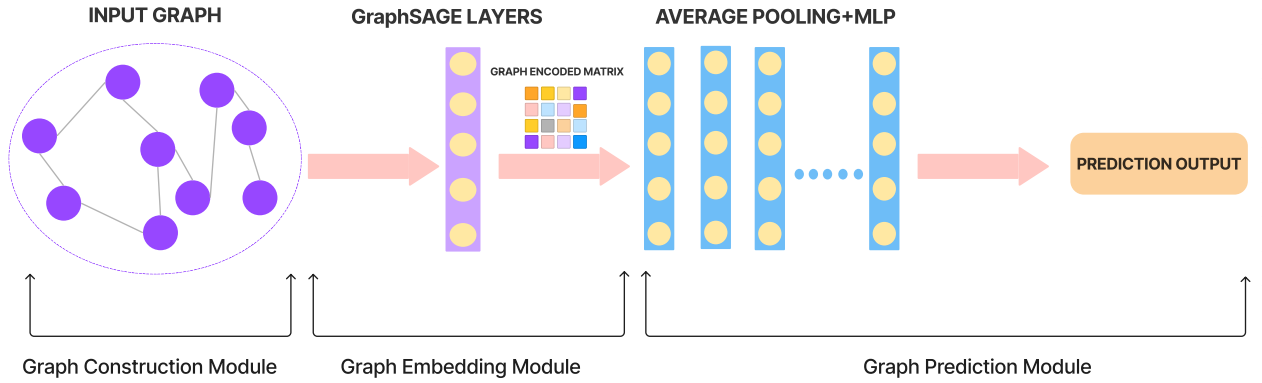


Fig. 3. The GNN framework for simile classification.

- (2) **The graph embedding module:** The authors of this research implemented Bi-Fuse GraphSAGE, a GNN framework for inductive representation learning of graphs which is used to generate low-dimensional vector representations for nodes. The embedding generation process takes the entire graph $G(V, E)$ and features for all nodes, $x_i \in V$ as input. In each iteration from $k = 0$ up to $k = K$ where k denotes the current step in the loop and $h_v^{(k)}$ denotes a node's representation at that step, K signifies the number of aggregator functions and W^k denotes set of the weight matrices in each iteration. First, each node $v \in V$ aggregates the representations of the nodes in its immediate neighbourhood, as represented by equation 1, into a single vector $h_{N(i)}^{(l+1)}$. After the aggregation of the neighbouring feature vectors, GraphSAGE concatenates the node's current representation h_v^k , with the aggregated neighbourhood vector $h_{N(v)}^{(k+1)}$, given in equation 2 and this concatenated

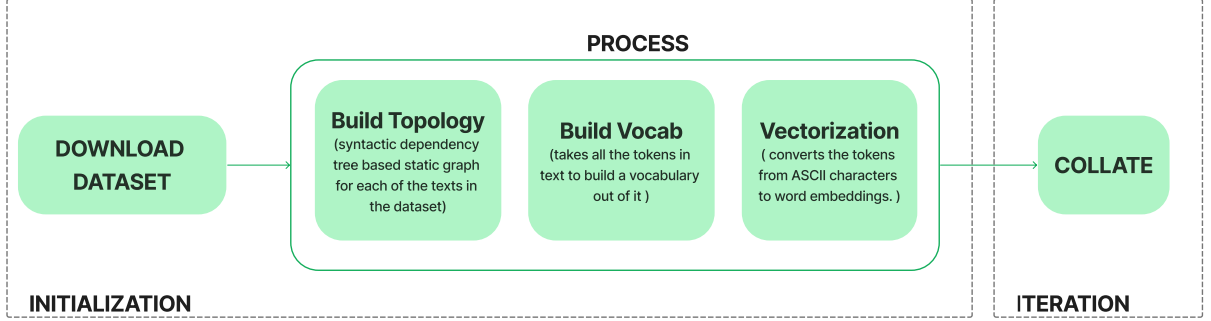


Fig. 4. Simile dataset preprocessing workflow.

vector is fed through a fully connected layer with a non-linear activation function represented by σ , following which each current node's representation is normalised as illustrated by equation 3.

$$h_{N(v)}^{(k+1)} = \text{aggregate}(h_v^k, \forall v \in N(v)) \quad (1)$$

$$h_v^{(k+1)} = \sigma(W^k \cdot \text{concat}(h_v^k, h_{N(v)}^{(k+1)} + b)) \quad (2)$$

$$h_v^{(k)} = \text{norm}(h_v^k) \quad (3)$$

- (3) **The prediction module:** The prediction module consists of an average pooling layer with 300 hidden units and an MLP classifier which produces predicted labels. In case of a presence of a simile is detected the framework predicts label '1' and in case of a non-simile, the framework predicts label '0'. The said trained framework is saved for predicting the presence of a simile on the sarcasm dataset. The training and validation accuracy and loss curves are discussed in section 5.

4.4.2. BERT-based structure for metaphor detection

Word2Vec and GloVe are two types of architectures used commonly for word embedding, but they have limitations in certain tasks such as representing terms that are not in their vocabulary and distinguishing between opposites. For example, the words "good" and "bad" are often very close to each other in the vector space created by these models, which can be problematic for NLP applications like sentiment analysis. On the other hand, BERT is a pretraining method that uses a self-supervised approach to learn from masked text sections. Developed by a team at Google Research, BERT is based on the transformer architecture and is designed to learn deep bidirectional representations from unlabeled text by conditioning on both left and right contexts. While Word2Vec and GloVe are unidirectional models that can only understand context in one direction, BERT can move sentences both to the left and right to fully comprehend the context of the target word or group of words. As a result, the detection of metaphors

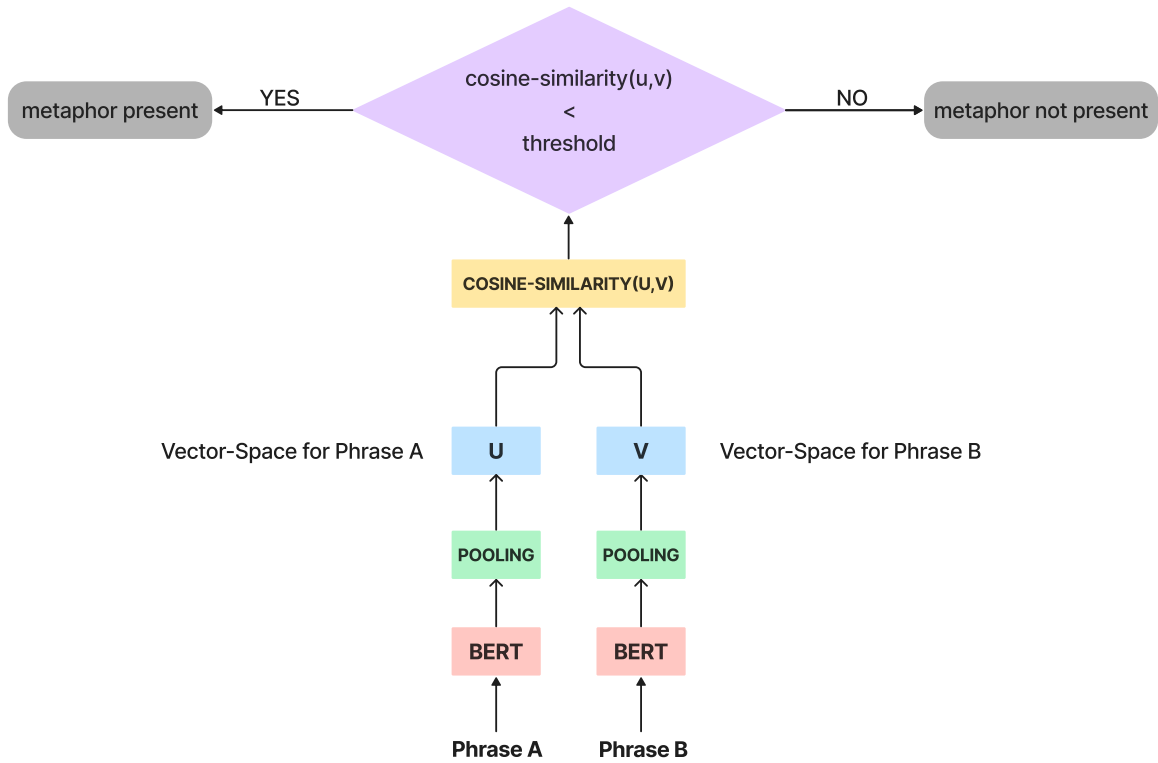


Fig. 5. The framework for metaphor detection.

is achieved by generating embeddings using a BERT-based network. The basic BERT model, without any fine-tuning, can generate embeddings for the phrases that can be used to calculate cosine similarity. Fig. 5 illustrates the framework used for detecting the presence of a metaphor and Fig. 6 illustrates the BERT-based network used for generating the embeddings with the hidden layer representations in red. For the BERT base, each encoder layer outputs a set of dense vectors.

Each vector contains 768 values each of which is nothing but contextual word embeddings. Initially, each sentence is split into two halves. For sentences of type I, phrase A consists of the subject and phrase B consists of the object. On the other hand, for sentences of type II, phrase A consists of a verb which acts on a noun phrase represented by phrase B. BERT-based embeddings are generated for both phrases A and B and cosine similarity is calculated. The entire process can be summarized as follows:

Detecting type-1 metaphors:

- (1) All the sentences in the dataset are first tokenized using spaCy.
- (2) The sentences are then split into two phrases: one containing the subject of comparison and the other containing the object of comparison.
- (3) Dense contextual embeddings are constructed for each of the phrases. The `last_hidden_state` tensor from the BERT model is extracted to quantify textual similarity. A pooling operation is performed that takes the mean of all token embeddings and compresses them into a single vector space representing a single phrase. Furthermore, the cosine similarity between the two vector spaces of the phrases is calculated. Following multiple trials, a threshold value of 0.7 was chosen to determine

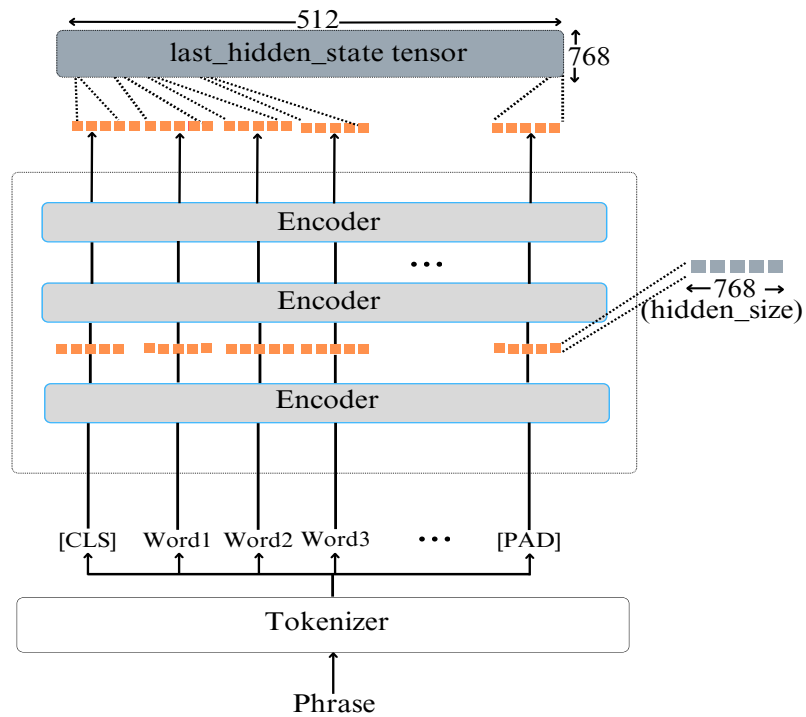


Fig. 6. BERT-based network used for generating embeddings.

the presence or absence of a metaphorical instance in a sentence. A cosine similarity larger than 0.7 accurately indicated the lack of a metaphorical statement, whereas one less than 0.7 suggested its presence.

Detecting type-2 metaphors:

- (1) All the sentences in the dataset are first tokenized using spaCy.
- (2) All the sentences are split into two phrases: one containing the personified verb and the other containing the object of comparison.
- (3) For each of the phrases, dense contextual embeddings are generated and textual similarity is measured in terms of cosine similarity.

4.4.3. Fuzzy logic-based approach for capturing polarity change in clauses

Fuzzy logic is a form of multi-valued logic that deals with reasoning that is approximate rather than precise. In fuzzy logic, a concept can possess a degree of truth (degree of membership) denoted by ' μ ' anywhere between 0.0 and 1.0, unlike in standard logic where a concept can only be completely true (1.0) or false (0.0). Fuzzy logic is useful for reasoning about inherently vague concepts, such as "tallness". In traditional set theory, an element either belongs to a set or not, but in fuzzy logic, an element can partially belong to a set, with a degree of membership ranging between 0 and 1. In general, the degree of membership is used to represent the uncertainty or fuzziness in a concept or variable and is a key component in fuzzy logic's ability to reason with imprecise or uncertain information. In fuzzy logic systems, a membership function is used to map this degree of membership of an element to a set,

based on a specified set of fuzzy rules. Trapezoidal membership functions are one way to define fuzzy sets. They are a type of membership function that allows for more flexibility in defining the shape of the fuzzy set.

A trapezoidal membership function is a piecewise linear function that has four parameters: a , b , c and d , where $a \leq b \leq c \leq d$. The function starts at 0 for $x \leq a$, increases linearly from 0 to 1 from a to b , stays at 1 from b to c , decreases linearly from 1 to 0 from c to d , and stays at 0 for $x \geq d$. By using trapezoidal membership functions, non-polygonal fuzzy sets can be defined that have smooth transitions between membership degrees. The shape of the membership function can be controlled by adjusting its parameters, allowing for the creation of fuzzy sets that match the intended intuition. For instance, a trapezoidal membership function can be used to define the fuzzy set "tall person" with parameters such as $a = 170$ cm, $b = 175$ cm, $c = 185$ cm and $d = 190$ cm. The resulting function would have a membership degree of 0 for heights below 170 cm and above 190 cm, a membership degree of 1 for heights between 175 cm and 185 cm, and smoothly increasing and decreasing membership degrees for heights between 170 cm and 175 cm and between 185 cm and 190 cm. Trapezoidal membership functions allow for the generation of non-polygonal fuzzy sets, which enable more precise and flexible modelling of linguistic variables in fuzzy logic.

The aim of this study is to assess the change in polarity from the main clause to the subordinate clause of a sentence, at different levels, in order to identify possible instances of sarcasm. Fuzziness is relevant in addressing the uncertainties related to the degree of positivity or negativity of an independent clause. By combining this with the computational identification of polarity shift, more effective outcomes can be achieved. The following steps were performed to detect sarcasm in the form of disparity of sentiments in clauses:

- (1) Tokenization: All the pre-processed textual data is first tokenised using spaCy's NLP object.
- (2) Separation of clauses: To find the polarity of a sentence's constituent clauses, a sentence is first separated into clauses. This is done in two ways. All the sentences are checked for the presence of subordinating conjunctions given in appendix A.3. If a sentence does not contain any of the subordinating conjunction mentioned then splitting of the sentence is done on the basis of two markers. The first marker is the subject of the sentence with syntactic dependency "nominal subject (nsubj)". The second marker is the last occurrence of any preposition in a sentence. It is marked with syntactic dependency "preposition (prep)". These markers divide the sentence into three parts, the first part spans from the beginning of the sentence to the first marker ("nsubj"), the second part spans from the first marker("nsubj") to the second marker("prep") and the third part spans from the second marker("prep") to the end of the sentence. For example, the sentence, "*You're everything I want in someone, I don't want anymore.*" splits into "*You're everything I want*" and "*in someone I don't want anymore*". Another example would be, "*Right before I die I am going to swallow a bag of popcorn kernels to make the cremation a bit more interesting.*" splits into "*Right before I die*", "*I am going to swallow a bag of popcorn kernels*" and "*to make the cremation a bit more interesting*".
- (3) Computing the polarity of clauses: For finding the polarity of a sentence's constituent clauses, pysentimiento [45] is used. pysentimiento is a python toolkit for sentiment analysis and text classification. It is a transformer-based open-source library. It uses BERTweet [46] as a base model in English. The list of constituent clauses is taken for each sentence and the polarity score corresponding to each clause in the form of positive, negative and neutral proportions is obtained.
- (4) Applying fuzzy logic to eliminate overlapping sentiment classes: In each sentence, every clause has three sentiment proportions, i.e., positive, negative and neutral. To determine the overall polarity of a clause fuzzy-logic has been implemented using Simpful [47]. Although pysentimiento provides

Table 5

The degree of each sentiment along with the corresponding polarity percentage associated with it

Sentiment	Degree	Range
Positive	Weakly positive	0-35 %
	Moderately positive	25-75 %
	Strongly positive	68-100 %
Negative	Weakly negative	2-34%
	Moderately negative	25-73 %
	Strongly negative	68-100 %
Neutral	Weakly neutral	1-35%
	Moderately neutral	25-73 %
	Strongly neutral	68-100 %

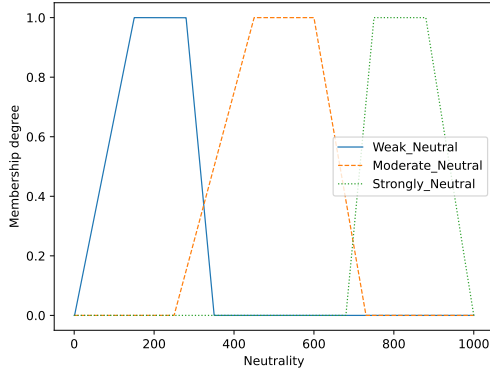
sentiment proportions for the positivity, negativity and neutrality of a clause, it does not provide any valuable information about the degree (weak, moderate, strong) of each sentiment. This study uses fuzzy logic to determine the overall polarity of the clauses with the help of a set of rules based on the projected degree of each sentiment. The degree of each sentiment namely, positive, negative and neutral have been devised using the assumed ranges given in table 5. The trapezoidal membership function is used to define a non-polygonal fuzzy set for each sentiment namely, positive, negative and neutral. Fig. 7 (a,b,c) illustrates various inputs for each sentiment.

The set of fuzzy rules used in this study can be found in appendix B.2. The fuzzified output is presented in (d) part of Fig. 7. Defuzzification is then applied to get the final polarity output.

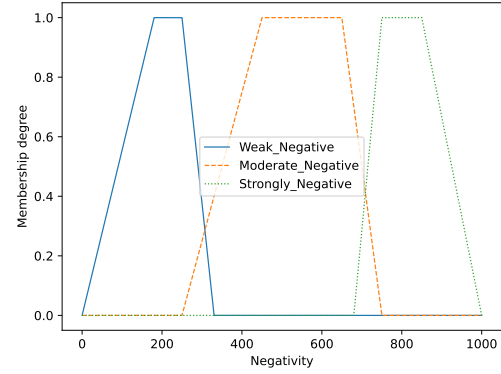
- (5) Checking for the disparity of sentiments: The total number of clauses with positive, neutral, or negative sentiment labels are counted and utilised to account for polarity shifts from negative to positive or positive to negative. If such a shift occurs, the function returns '1'; otherwise, it returns '0'.

4.5. DWAEF: The deep weighted average ensemble framework

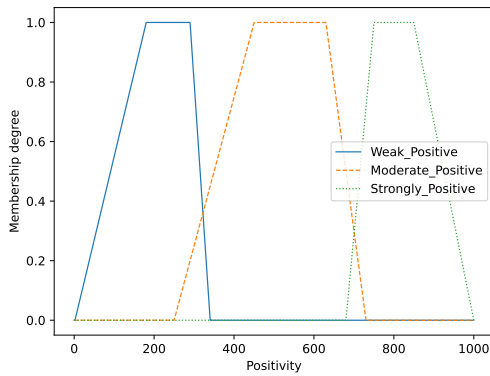
DWAEF, the proposed deep weighted average ensemble-based framework, is a three-tiered structure. It comprises three base learners: a TabNet, a 1-D CNN and an MLP. The framework is implemented with the help of DeepStack Library [48]. Initially, the curated dataset is used to pretrain the three models. During training, each of the base learners receives the outputs of the GNN-based simile detection framework, the BERT-based metaphor detection framework, the fuzzy-based polarity shift detection framework as well as the primitive features as its inputs. Next, the predictions produced by each module of the ensemble are weighed based on their performance on a hold-out validation dataset. This is represented by a Dirichlet ensemble object. To achieve more accurate solution a weight optimization search is used in conjunction with randomized search based on the Dirichlet distribution on the validation dataset. Randomized search based on the Dirichlet distribution involves randomly sampling the weights from a Dirichlet distribution, which generates a probability distribution over the weights that is continuous and flexible. Weight optimization search iteratively adjusts the ensemble model weights to maximize accuracy on the validation dataset. Through this process the base learners, a TabNet, a 1-D CNN and an MLP, are assigned optimal weights. These weights determine the contribution of each model to the final prediction. Subsequently, the predictions of each model are combined through a weighted average. No meta-learner is used in this ensemble method. Lastly, the ensemble model is assessed on a separate test set to determine how well it can perform on new data.



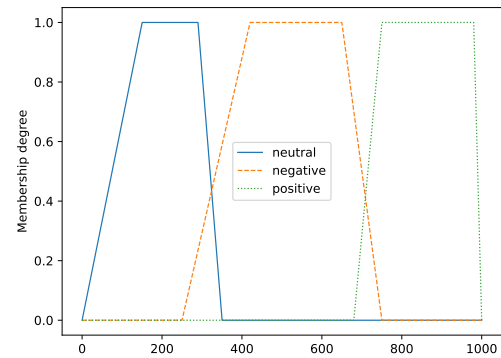
(a) The inputs for the neutral sentiment



(b) The inputs for the negative sentiment



(c) The inputs for the positive sentiment



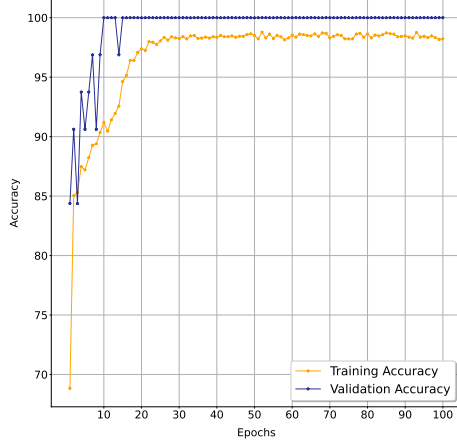
(d) The fuzzified output

Fig. 7. Fuzzy inputs and output for neutral, negative and positive sentiments.

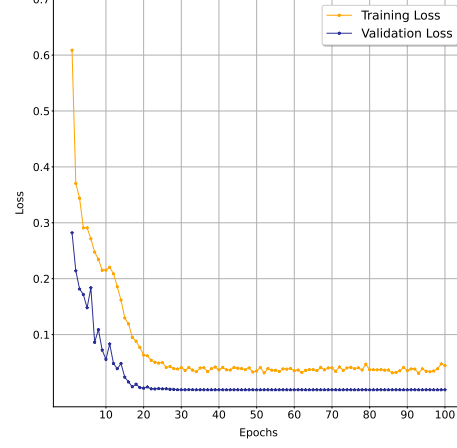
The findings that were achieved through the utilization of the suggested methodology are reported in section 5.

5. Evaluation, results and analysis

This section examines the results acquired by employing different techniques on the dataset. The purpose of the proposed methodology is to efficiently detect sarcasm in online text using the presence of figurative comparisons, i.e., similes and metaphors and shifts in polarity of the text's constituent clauses. This segment is organised as follows: Subsection 5.1 discusses the accuracy vs. epochs and loss vs. epoch curves for the GNN framework. The accuracy scores obtained during experimentation with different threshold values by the BERT-based Metaphor Detection Framework are discussed in subsection 5.2. The confusion metrics for the fuzzy-based approach are presented in subsection 5.3. Subsection 5.4 analyses the results obtained using DWAEF, the deep weighted average ensemble model.



(a) Training accuracy (in orange) and validation accuracy (in purple) vs. number of epochs



(b) Training loss(in orange) and validation loss (in purple) vs. number of epochs

Fig. 8. Training and validation curves of the GNN framework for simile detection.

Table 6
Hyperparameter settings for the GNN framework

Parameter	Value
seed	1234
batch size	32
epochs	100
learning rate	0.01
learning rate patience	2
learning rate reduce factor	0.5
hidden layers	300
dropout	0.3
graph pooling	Average Pooling
optimizer	Adam
loss	Cross Entropy
activation	ReLU

5.1. Results obtained by the GNN-based simile detection Framework

The GNN framework described in the section 4.4.1 was pretrained on a dataset of roughly 3,000 sentences, out of which roughly 50% were similes and the rest 50% were non-similes. The entire dataset was split into three groups: 60% of the data was used for training, 20% was used for testing and the rest 20% was used for validation. This pretrained framework was then deployed on the main collated sarcasm dataset to extract the presence of a simile as one of the features. With a batch size of 32, the model was executed for 100 epochs. The rest of the hyperparameters are given in table 6. The training and validation curves of the proposed framework are illustrated in Fig. 8. It is evident from both curves that the framework is free from both overfitting and underfitting.

The testing accuracy obtained using the proposed GNN framework for simile detection was 99.22% using GloVe word embeddings. The state-of-the-art results ensured accurate detection of the simile in

Table 7
Various threshold values and their corresponding accuracy score

Threshold value	Accuracy	
	Type1	Type2
0.3	0.72	0.60
0.4	0.62	0.60
0.5	0.70	0.68
0.6	0.78	0.66
0.7	0.82	0.78

the main sarcasm dataset.

5.2. Results obtained by the BERT-based metaphor detection Framework

Before settling on the best threshold value to assess the existence or absence of a metaphorical instance in a sentence, several values were tested. The various values tested and the accompanying accuracy values are listed in table 7. It is clear that at a threshold value of 0.7, the most accurate predictions were achieved for both type 1 and type 2 metaphors. Thus, a cosine similarity of more than 0.7 accurately indicates the absence of a metaphorical remark, whereas a cosine similarity of less than 0.7 indicates its presence.

5.3. Results obtained by the fuzzy-based polarity shift detection Framework

The confusion matrix shown in Fig. 9 depicts the performance of the fuzzy-based framework. In the matrix, there are four distinct combinations of expected and actual values. It is evident from the matrix that the framework properly identified the variations in polarity that actually indicated sarcasm at the clausal structural level of 1,309 sentences. Additionally, 125 incorrect sentences were identified as true. On the other hand, it also correctly identified the lack of a tone change in 1,228 sentences but failed to recognise a polarity shift in 127 sentences. It is obvious from the matrix that the framework made accurate predictions for 87.81% of the whole dataset. While 91.28% of all predicted true classes were predicted actually true, 85.22% of all real true classes were predicted true by the framework. One of the reasons for the state-of-the-art outcome is the usage of fuzzy rules to cope with uncertainties over whether a solitary sentence is positive or negative.

5.4. Comparison of results obtained by individual base learners and DWAEF with primitive and proposed features

The study performs and compares the results of three experiments- 1: sarcasm detection using individual base learners, namely- TabNet, 1-D CNN and MLP with primitive features; 2: sarcasm detecton using the same base learners with a combination of both primitive and proposed features; and 3: sarcasm detection using DWAEF with a combination of primitive and proposed features. The hyperparameter settings of the base learners in all the three experiments were kept the same to ensure that any observed differences in performance were due to the changes in the feature set rather than differences in hyperparameters. Table 8 gives the corresponding hyperparameter settings for each of the base learners.

The results of experiment 1 and 2 are displayed in table 9. The comparison shows that the outcome of experiment 2 is better than that of experiment 1. Hence, the accuracy values of experiment 2 is then compared with the accuracy value of the third experiment (the overall accuracy of DWAEF) and also

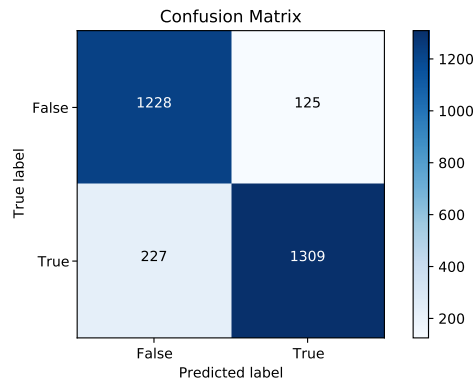


Fig. 9. Confusion matrix for the fuzzy-based polarity shift detection framework

with the three most powerful and widely used traditional machine learning classifiers, RF, SVM and AdaBoost (AB). The corresponding detailed comparison report is summarised in table 10.

A confusion matrix depicting the performance of DWAEF is presented in Fig 10. Based on the presented confusion matrix, the performance of the model can be inferred as follows. The model achieved a high accuracy rate of 92.01%, indicating that it correctly predicted the majority of cases. Additionally, the precision of the model was found to be 92.98%, indicating that when it predicted a positive class, it was correct 92.98% of the time. Furthermore, the recall rate of the model was 89.59%, indicating that the model correctly identified 89.59% of actual positive cases. The F1 score of the model was also found to be high at 91.62%, indicating that the model performed well in terms of both precision and recall. However, the model exhibited some limitations, such as a number of false positives (FP) and false negatives (FN). Specifically, the model produced 91 false positive predictions, which could lead to unnecessary actions or decisions being taken. Additionally, the model produced 140 false negative predictions, which could result in missed opportunities or errors in decision-making. In conclusion, while the model performed well with high accuracy, precision, recall and F1 score, there is still room for improvement in reducing the number of false positives and false negatives to further enhance its performance. These findings may have implications for decision-making processes in applications that utilize this model, and can inform future improvements in its development.

It can be inferred that the DWAEF outperforms all the models used in this study in terms of accuracy. Also, the proposed features of this research, namely, presence of figurative comparisons, i.e., simile and metaphor and shift in polarity of a sentence's constituent clauses, successfully aid in better detection of sarcasm in online text messages. The detection of sarcasm has also been boosted by switching to a deep weighted average ensemble framework because the framework assigns each base member's share of the prediction weight based on how well it performed individually during training. Moreover, researchers in [22] failed to detect sarcasm in sentences written in a formal and polite tone. However, including the proposed novel indicators successfully detected sarcasm in such sentences. Table 11 presents some of them.

6. Conclusion and future work

Detection of sarcasm poses one of the leading challenges in sentiment analysis, as even a single sarcastic remark can influence sentiment analyzers to produce undesirable results. Primitive techniques used

Table 8
Hyperparameter settings for TabNet, 1-D CNN and MLP used in DWAEF

Model	Hyperparameter settings
TabNet	optimizer: Adam learning rate: 0.001 step size:10 gamma:1.4 mask type: entmax
1-D CNN	seed:1234 learning rate: 0.0025 dropout rate: 0.8 loss: sparse categorical cross entropy optimizer: SGD
MLP	activation: ReLU alpha: 0.00025 hidden layer sizes: (200,150,100,50,25) learning rate: Adaptive solver: Adam maximum iteration: 200 random state: 25

Table 9
Accuracy scores for DWAEF base learners- TabNet, 1-D CNN and MLP

Stage	TabNet	1-D CNN	MLP
Stage 1: Primitive features only	76.80	75.03	64.09
Stage 2: Primitive features + proposed features	89.80	87.07	85.21

Table 10
Summary of accuracy scores of all classifiers

Model	Accuracy score(%)
RF	81.37
SVM	78.58
AB	81.13
MLP	85.21
TabNet	89.80
1-D CNN	87.07
DWAEF	92.01

Table 11
Examples of correctly classified sentences written in a polite/formal tone

Sentences	Novel indicator present	Classification result(%)
Everyone has a photographic memory, some just don't have film.	Metaphor	1
When it comes to finding a good place to eat you can't doubt her choice, she's a connoisseur of food no wonder why she eats like a pig.	Simile, Clauses	1
No, you're right, we should just put the mentally ill down like dogs if they do something inappropriate.	Simile	1

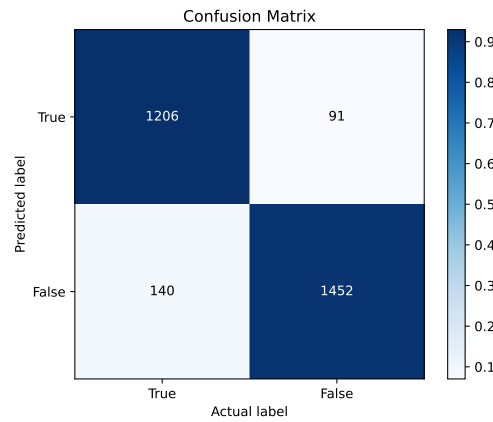


Fig. 10. Confusion matrix for DWAEF

in sarcasm detection used low-level features and traditional machine learning algorithms. The study looked into sarcasm detection with a new perspective. It proposed DWAEF, a deep-weighted ensemble-based framework for sarcasm detection. The framework utilized figurative speech components mainly, the presence of simile, the presence of metaphor and the change in the polarity of a sentence's constituent clauses using deep learning techniques. The predictions done by the above modules were then fed into DWAEF, which comprised a 1-D CNN, a TabNet and an MLP as its base learners. Based on the results, it can be concluded that combining the proposed indicators with the primitive features achieved better results across all classifiers. It was seen that the proposed ensemble framework performed better as compared to traditional machine learning classifiers. The proposed technique achieved the highest accuracy of 92.01% when proposed indicators were combined with primitive features and evaluated using a weighted average ensemble of deep learning algorithms. The study employed various state-of-the-art tools and techniques; still, the proposed framework may be made to improve the model's characteristics and efficiency.

In future, the researchers plan to investigate the use of larger datasets for sarcasm detection in deep ensemble framework, while considering potential drawbacks such as increased time and complexity for training and a lack of annotation resources. Efficient resource management strategies such as unsupervised learning techniques and pre-trained models, and crowd sourcing model for annotation shall be explored to mitigate these challenges. The optimal dataset size will depend on the specific application and available resources, and further research should identify the trade-offs and limitations associated with using larger datasets for sarcasm detection.

While working in the domain of sarcasm detection, the authors also came across other crucial perspectives regarding sarcasm. It was sensed that the interpretation of sarcasm can also be influenced by cultural and linguistic factors, and people from different cultures may interpret sarcasm differently. Different languages may have their own unique styles and patterns of sarcasm, which may not be easily translatable or understandable to those unfamiliar with the language. Therefore, cultural and linguistic contexts are crucial when studying or detecting sarcasm. Sarcasm detection models that work well in one language or culture may not perform well in another, highlighting the importance of considering these factors in model design and evaluation. Hence, in future, the authors plan to incorporate more advanced

tools in the DWAEF framework and fine-tune it to perform cross-lingual, cross-cultural and multimodal predictions.

Appendices

A. List of intensifiers, interjections and subordinating conjunction

A.1. List of intensifiers

- | | |
|-----------------|-------------------|
| ● amazingly | ● extremely |
| ● ass | ● extraordinarily |
| ● astoundingly | ● fantastically |
| ● awful | ● frightfully |
| ● bare | ● fucking |
| ● bloody | ● fully |
| ● crazy | ● hella |
| ● dead | ● incredibly |
| ● dreadfully | ● insanely |
| ● colossally | ● literally |
| ● especially | ● mad |
| ● exceptionally | ● mightily |
| ● excessively | ● most |

A.2. List of Interjections

- | | |
|-------------------|----------------|
| ● aha | ● darn |
| ● ahem | ● duh |
| ● ahh | ● eek |
| ● ahoy | ● eh |
| ● alas | ● encore |
| ● arg | ● eureka |
| ● aw | ● fiddlesticks |
| ● bam | ● gadzooks |
| ● bingo | ● gee |
| ● blah | ● gee |
| ● boo | ● whiz |
| ● bravo | ● golly |
| ● brrr | ● goodbye |
| ● cheers | ● goodness |
| ● congratulations | ● good grief |
| ● dang | ● gosh |
| ● drat | ● ha-ha |

- 1 ● hallelujah
- 2 ● hello
- 3 ● hey
- 4 ● hmm
- 5 ● holy buckets
- 6 ● holy cow
- 7 ● holy smokes
- 8 ● hot dog
- 9 ● huh
- 10 ● humph
- 11 ● hurray

- 1 ● oh
- 2 ● oh dear
- 3 ● oh my
- 4 ● oh well
- 5 ● oops
- 6 ● ouch
- 7 ● ow
- 8 ● phew
- 9 ● phooey
- 10 ● pooh
- 11 ● pow

12 A.3. List of subordinating conjunctions

- 15 ● after
- 16 ● before
- 17 ● as soon as
- 18 ● while
- 19 ● when
- 20 ● as
- 21 ● because
- 22 ● since
- 23 ● if
- 24 ● provided that
- 25 ● as long as
- 26 ● unless

- 15 ● although
- 16 ● though
- 17 ● even though
- 18 ● then
- 19 ● which
- 20 ● who
- 21 ● that
- 22 ● whose
- 23 ● and
- 24 ● but
- 25 ● &

28 B. Implementation details

29 B.1. The initialization preprocessing steps of GNN-based simile detection framework

30 For pretraining the GNN-based framework for simile detection the present study curated a dataset
 31 comprising approximately 3,512 English language sentences that have been systematically categorized
 32 into two groups, i.e., sentences containing similes and those that do not. The dataset can be accessed
 33 from <https://github.com/simdeol/Simile-Dataset>. The dataset was curated through online sources and it
 34 underwent a rigorous process of double-annotation to ensure its accuracy and reliability. This dataset
 35 offers the opportunity for future research on the impact of similes in both written and spoken language,
 36 thereby contributing to a better understanding of figurative language in general. Following preprocess-
 37 ing, 3,000 sentences were selected for further analysis. The preprocessing steps illustrated in Fig. 4 are
 38 described below:

- 39 ● `build_topology`: This step builds a syntactic dependency text graph for each of the data items in the
 40 raw dataset.
- 41 ● `build_vocab`: This step is responsible for building vocabulary out of all tokens appearing in the data
 42 items.

- vectorization: This step is the lookup step, responsible for converting tokens from ASCII characters to word embeddings

B.2. The fuzzy rule set for detecting polarity shift in clauses

- (1) IF (Neutral IS Weak_Neutral) AND (Negative IS Weak_Negative) AND (Positive IS Moderate_Positive) THEN (Output IS positive)
- (2) IF (Neutral IS Weak_Neutral) AND (Positive IS Weak_Positive) AND (Negative IS Moderate_Negative) THEN (Output IS negative)
- (3) IF (Positive IS Weak_Positive) AND (Negative IS Weak_Negative) AND (Neutral IS Moderate_Neutral) THEN (Output IS neutral)
- (4) IF (Positive IS Weak_Positive) AND (Negative IS Moderate_Negative) AND (Neutral IS Moderate_Neutral) THEN (Output IS negative)
- (5) IF (Negative IS Weak_Negative) AND (Positive IS Moderate_Positive) AND (Neutral IS Moderate_Neutral) THEN (Output IS positive)
- (6) IF (Neutral IS Weak_Neutral) AND (Negative IS Moderate_Negative) AND (Positive IS Moderate_Positive) THEN (Output IS positive)
- (7) IF (Neutral IS Strongly_Neutral) AND (Negative IS Weak_Negative) AND (Positive IS Weak_Positive) THEN (Output IS neutral)
- (8) IF (Positive IS Strongly_Positive) AND (Negative IS Weak_Negative) AND (Neutral IS Weak_Neutral) THEN (Output IS positive)
- (9) IF (Negative IS Strongly_Negative) AND (Neutral IS Weak_Neutral) AND (Positive IS Weak_Positive) THEN (Output IS negative)

References

- [1] E.W. Pamungkas, V. Basile and V. Patti, A joint learning approach with knowledge injection for zero-shot cross-lingual hate speech detection, *Inf. Process. Manag.* **58**(4) (2021), 102544. doi:10.1016/j.ipm.2021.102544.
- [2] S. Frenda, A.T. Cignarella, V. Basile, C. Bosco, V. Patti and P. Rosso, The unbearable hurtfulness of sarcasm, *Expert Systems with Applications* **193** (2022), 116398.
- [3] V. Basile, It's the end of the gold standard as we know it. On the impact of pre-aggregation on the evaluation of highly subjective tasks, in: *Proceedings of the AIXIA 2020 Discussion Papers Workshop co-located with the 19th International Conference of the Italian Association for Artificial Intelligence (AIXIA2020)*, Vol. 2776, CEUR-WS, 2020, pp. 31–40. <http://ceur-ws.org/Vol-2776/paper-4.pdf>.
- [4] M. Gori, G. Monfardini and F. Scarselli, A new model for learning in graph domains, in: *Proceedings. 2005 IEEE International Joint Conference on Neural Networks*, Vol. 2, 2005, pp. 729–734.
- [5] J. Devlin, M.-W. Chang, K. Lee and K. Toutanova, BERT: pre-training of deep bidirectional transformers for language understanding, *arXiv preprint arXiv:1810.04805* (2018).
- [6] L.A. Zadeh, Fuzzy logic, *Computer* **21**(4) (1988), 83–93.
- [7] S.Ö. Arik and T. Pfister, Tabnet: attentive interpretable tabular learning, in: *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 35, 2021, pp. 6679–6687.
- [8] A. Reyes and P. Rosso, Making objective decisions from subjective data: detecting irony in customer reviews, *Decision Support Systems* **53**(4) (2012), 754–760. <https://doi.org/10.1016/j.dss.2012.05.027>.
- [9] S. Attardo, Sarcasm and similes, *Journal of Pragmatics* **27**(5) (1997), 683–700.
- [10] E. Riloff, A. Qadir, P. Surve, L. De Silva, N. Gilbert and R. Huang, Sarcasm as contrast between a positive sentiment and negative situation, in: *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, 2013, pp. 704–714.
- [11] S.M. Sarsam, H. Al-Samarraie, A.I. Alzahrani and B. Wright, Sarcasm detection using machine learning algorithms in Twitter: a systematic review, *International Journal of Market Research* **62**(5) (2020), 578–598.

- [12] M. Bouazizi and T.O. Ohtsuki, A pattern-based approach for sarcasm detection on Twitter, *IEEE Access* **4** (2016), 5477–5488.
- [13] K. Hallmann, F. Kunneman, C. Liebrecht, A. van den Bosch and M. van Mulken, Sarcastic Soulmates: Intimacy and irony markers in social media messaging, in: *Linguistic Issues in Language Technology, Volume 14, 2016-Modality: Logic, Semantics, Annotation, and Machine Learning*, 2016.
- [14] N. Littlestone, Learning quickly when irrelevant attributes abound: a new linear-threshold algorithm, *Machine Learning* **2**(4) (1988), 285–318.
- [15] D. Bamman and N. Smith, Contextualized sarcasm detection on Twitter, in: *Proceedings of the International AAAI Conference on Web and Social Media*, Vol. 9, 2015, pp. 574–577.
- [16] S. Lukin and M. Walker, Really? well. apparently bootstrapping improves the performance of sarcasm and nastiness classifiers for online dialogue, *arXiv preprint arXiv:1708.08572* (2017).
- [17] S. Muresan, R. Gonzalez-Ibanez, D. Ghosh and N. Wacholder, Identification of nonliteral language in social media: a case study on sarcasm, *Journal of the Association for Information Science and Technology* **67**(11) (2016), 2725–2737.
- [18] P. Mehndiratta and D. Soni, Identification of sarcasm using word embeddings and hyperparameters tuning, *Journal of Discrete Mathematical Sciences and Cryptography* **22**(4) (2019), 465–489.
- [19] M.S. Razali, A.A. Halin, L. Ye, S. Doraisamy and N.M. Norowi, Sarcasm detection using deep learning with contextual features, *IEEE Access* **9** (2021), 68609–68618.
- [20] A. Baruah, K. Das, F. Barbhuiya and K. Dey, Context-aware sarcasm detection using BERT, in: *Proceedings of the Second Workshop on Figurative Language Processing*, 2020, pp. 83–87.
- [21] M. Abdullah, J. Khrais and S. Swedat, Transformer-based deep learning for sarcasm detection with imbalanced dataset: resampling techniques with downsampling and augmentation, in: *2022 13th International Conference on Information and Communication Systems (ICICS)*, IEEE, 2022, pp. 294–300.
- [22] P. Goel, R. Jain, A. Nayyar, S. Singhal and M. Srivastava, Sarcasm detection using deep learning and ensemble learning, *Multimedia Tools and Applications* (2022), 1–24.
- [23] J. Lemmens, B. Burtenshaw, E. Lotfi, I. Markov and W. Daelemans, Sarcasm detection using an ensemble approach, in: *Proceedings of the Second Workshop on Figurative Language Processing*, 2020, pp. 264–269.
- [24] J. Plepi and L. Flek, Perceived and intended sarcasm detection with graph attention networks, *arXiv preprint arXiv:2110.04001* (2021).
- [25] L. Wu, Y. Chen, K. Shen, X. Guo, H. Gao, S. Li, J. Pei and B. Long, Graph neural networks for natural language processing: a survey, *arXiv preprint arXiv:2106.06090* (2021).
- [26] S. He, F. Guo and S. Qin, Sarcasm detection using graph convolutional networks with bidirectional LSTM, in: *Proceedings of the 2020 3rd International Conference on Big Data Technologies*, 2020, pp. 97–101.
- [27] C. Lou, B. Liang, L. Gui, Y. He, Y. Dang and R. Xu, Affective dependency graph for sarcasm detection, in: *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*, 2021, pp. 1844–1849.
- [28] P. Chaudhari and C. Chandankhede, Literature survey of sarcasm detection, in: *2017 International Conference on Wireless Communications, Signal Processing and Networking (WiSPNET)*, IEEE, 2017, pp. 2041–2046.
- [29] D.G. Maynard and M.A. Greenwood, Who cares about sarcastic tweets? Investigating the impact of sarcasm on sentiment analysis, in: *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC'14)*, ELRA, 2014.
- [30] A. Rajadesingan, R. Zafarani and H. Liu, Sarcasm detection on Twitter: a behavioral modeling approach, in: *Proceedings of the Eighth ACM International Conference on Web Search and Data Mining*, 2015, pp. 97–106.
- [31] F. Barbieri, F. Ronzano and H. Saggion, UPF-talm: SemEval 2015 tasks 10 and 11. Sentiment analysis of literal and figurative language in Twitter, in: *Proceedings of the 9th International Workshop on Semantic Evaluation (SemEval 2015)*, 2015, pp. 704–708.
- [32] V. Niculae and C. Danescu-Niculescu-Mizil, Brighter than gold: figurative language in user generated comparisons, in: *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2014, pp. 2008–2018.
- [33] J. O'Donoghue, Is a metaphor (like) a simile? Differences in meaning, effect and processing, *UCL Working Papers in Linguistics* **21** (2009), 125–149.
- [34] M. Popa-Wyatt, Ironic metaphor: a case for metaphor's contribution to truth-conditions, in: *E. Walaszewska, M. Kisielewska-Krysiuk & A. Piskorska (ed.) In the Mind and Across Minds: A Relevance-theoretic Perspective on Communication and Translation*, 2010.
- [35] S. Krishnakumaran and X. Zhu, Hunting elusive metaphors using lexical resources., in: *Proceedings of the Workshop on Computational Approaches to Figurative Language*, 2007, pp. 13–20.
- [36] X. Dong, Z. Yu, W. Cao, Y. Shi and Q. Ma, A survey on ensemble learning, *Frontiers of Computer Science* **14**(2) (2020), 241–258.
- [37] S. Loria, TextBlob documentation, *Release 0.15* **2** (2018).

- 1 [38] R. Misra and P. Arora, Sarcasm detection using hybrid neural network, *arXiv preprint arXiv:1908.07414* (2019). 1
- 2 [39] M. Khodak, N. Saunshi and K. Vodrahalli, A large self-annotated corpus for sarcasm, *arXiv preprint arXiv:1704.05579* 2
(2017).
- 3 [40] J. Eisenstein, What to do about bad language on the internet, in: *Proceedings of the 2013 Conference of the North American* 3
Chapter of the Association for Computational Linguistics: Human Language Technologies, Association for Computational
4 Linguistics, Atlanta, Georgia, 2013, pp. 359–369. <https://aclanthology.org/N13-1037>. 4
- 5 [41] C. Liebrecht, F. Kunneman and A. van Den Bosch, The perfect solution for detecting sarcasm in tweets# not (2013). 5
- 6 [42] C.D. Manning, M. Surdeanu, J. Bauer, J.R. Finkel, S. Bethard and D. McClosky, The Stanford CoreNLP natural language 6
processing toolkit, in: *Proceedings of 52nd Annual Meeting of the Association for Computational Linguistics: System*
7 *Demonstrations*, 2014, pp. 55–60. 7
- 8 [43] W. Hamilton, Z. Ying and J. Leskovec, Inductive representation learning on large graphs, *Advances in Neural Information* 8
Processing Systems **30** (2017). 9
- 10 [44] L. Wu, Y. Chen, H. Ji and B. Liu, Deep learning on graphs for natural language processing, in: *Proceedings of the 44th* 10
International ACM SIGIR Conference on Research and Development in Information Retrieval, 2021, pp. 2651–2653. 11
- 11 [45] J.M. Pérez, J.C. Giudici and F. Luque, pysentimiento: a Python toolkit for sentiment analysis and socialNLP tasks (2021). 11
- 12 [46] D.Q. Nguyen, T. Vu and A.T. Nguyen, BERTweet: A pre-trained language model for English tweets, *arXiv preprint* 12
arXiv:2005.10200 (2020). 13
- 13 [47] S. Spolaor, C. Fuchs, P. Cazzaniga, U. Kaymak, D. Besozzi and M.S. Nobile, Simpful: a user-friendly Python library for 13
fuzzy logic, *International Journal of Computational Intelligence Systems* **13**(1) (2020), 1687–1698. 14
- 14 [48] J. Borges, DeepStack: ensembles for deep learning, 2019. <https://github.com/jcborges/DeepStack>. 14
- 15 15
- 16 16
- 17 17
- 18 18
- 19 19
- 20 20
- 21 21
- 22 22
- 23 23
- 24 24
- 25 25
- 26 26
- 27 27
- 28 28
- 29 29
- 30 30
- 31 31
- 32 32
- 33 33
- 34 34
- 35 35
- 36 36
- 37 37
- 38 38
- 39 39
- 40 40
- 41 41
- 42 42
- 43 43
- 44 44
- 45 45
- 46 46