

# A Multi-type Classifier Ensemble for Detecting Fake Reviews Through Textual-based Feature Extraction

GREGORIUS SATIA BUDHI

School of Information and Physical Sciences, The University of Newcastle, Callaghan, NSW 2308, Australia,

[Gregorius.Satiabudhi@uon.edu.au](mailto:Gregorius.Satiabudhi@uon.edu.au)

Informatics Department, Petra Christian University, Surabaya 60236, Indonesia, [Greg@petra.ac.id](mailto:Greg@petra.ac.id)

RAYMOND CHIONG\*

School of Information and Physical Sciences, The University of Newcastle, Callaghan, NSW 2308, Australia,

[Raymond.Chiong@newcastle.edu.au](mailto:Raymond.Chiong@newcastle.edu.au)

The financial impact of online reviews has prompted some fraudulent sellers to generate fake consumer reviews for either promoting their products or discrediting competing products. In this study, we propose a novel ensemble model—the Multi-type Classifier Ensemble (MtCE)—combined with a textual-based featuring method, which is relatively independent of the system, to detect fake online consumer reviews. Unlike other ensemble models that utilise only the same type of single classifier, our proposed ensemble utilises several customised machine learning classifiers (including deep learning models) as its base classifiers. The results of our experiments show that the MtCE can adequately detect fake reviews, and that it outperforms other single and ensemble methods in terms of accuracy and other measurements in all the relevant public datasets used in this study. Moreover, if set correctly, the parameters of MtCE, such as base-classifier types, the total number of base classifiers, bootstrap and the method to vote on output (e.g., majority or priority), further improve the performance of the proposed ensemble.

**Keywords:** Fake review detection, online commerce security, novel ensemble model, machine learning, deep learning.

## 1 INTRODUCTION

Fake review detection is a hot research topic in natural language processing [55]. A fake review is a consumer review of a product with fictitious opinions written deliberately to sound authentic, for commercial motives; that is, to promote or damage the reputation of the reviewed product [42; 50]. There is a significant possibility that fake reviews distort the actual evaluation of a product [23], create harmful effects and eventually reduce trust in consumer reviews and undermine the effectiveness of the online market [43]. These fraudulent acts occur in response to the vital role of consumer reviews in purchasing behaviour in online markets [68]; consumers trust opinions expressed in online reviews and depend on them to make decisions [4; 14; 39]. Without detection efforts, reviews may be replete with lies, fakes and deceptions, and thus completely useless—or worse [55].

The majority of studies in fake review detection follow three main approaches: those based on the content of the reviews, those based on the behaviour of the reviewers, or a combination of these [3; 29]. Content-based methods focus on the linguistic features of text such as words, parts of speech (POS), n-gram, term frequency (TF) or other linguistic characteristics [21; 27; 55]. The behaviour-based approach focuses more on reviewers' behaviour, such as the user's identity, reviewed product, total number of reviews, ratings given, and duration of reviews [1; 3; 59]. This approach is straightforward, needs only a small number of features, and can deliver good results if properly designed. However, the behaviour-based approach is entirely dependent on additional information, such as metadata, provided by the system. In contrast, the content-based approach is independent of the system, requiring only the review text.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

© 2022 Association for Computing Machinery.

1533-5399/2022/1-ART1 \$15.00

<http://dx.doi.org/10.1145/3568676>

The content-based approach has two sub-categories. The first creates input features from the review text using Bag of Words (BOW), Word2Vec and skip-gram; thus, the input features for detection are words or terms created from the review text itself, and we call this approach textual-based featuring. The features extracted by this approach are more or less similar, mainly in the form of n-gram terms [14; 21; 27; 42; 63]. This textual-based approach is promising and is used in many studies to extract features from review texts. While independent of the system and needing only the text, this approach usually produces good results [11; 21; 27; 41; 42; 49; 55; 63; 72]. The second form of content-based method attempts to extract information and properties from the text, such as the length of the text, total words and total sentences, in addition to linguistic characteristics of the text, such as POS, subjectivity, complexity, diversity and similarity, and sentiment analysis [12; 26; 28; 57; 65; 66; 71]. This type of content-based featuring is lightweight compared with textual-based featuring; for example, only 80 features in [12] compared with 5,000 features in [11]. This approach, too, is independent of the system since it requires only the text; however, performance is usually relatively poor if not combined with other approaches [12].

In this study, we focus on the textual-based approach. For our experiments, we used Ott et al.'s [49] and Li et al.'s [41] datasets. These datasets are public fake review datasets used in many studies (e.g., see [12; 21; 27; 42; 55; 57; 72]). Before extracting the input features, we implemented text preprocessing such as tokenisation, stopword removal, detection of negation words, correction of elongation words, and POS lemmatisation. To extract the input features, we used the BOW method, n-gram words (from unigram to trigram words) and the count vectorised method. These methods convert a collection of text documents to a matrix of token counts. The token count features are in descending order by TF across the corpus (dataset). To detect fake reviews, we propose a novel machine learning (ML) ensemble model—the Multi-type Classifier Ensemble (MtCE)—that utilises several different types of standard (single) classifiers as its base classifier.

Ensemble classifiers are increasingly used in many different areas of study, such as text analysis [10; 39], computer security [30; 31], environmental science [62], medical research [69], electrical fault detection [47] and stock market prediction [34]. The main idea behind ensemble classifier models is to combine multiple single classifiers to produce a combination that outperforms any single classifier itself [8; 56]. In general, ensemble classifier models can be classified into four groups: bagging, boosting, stacked generalisation, and the mixture of experts [52]. Bagging, also called bootstrap aggregating, accumulates multiple predictors/classifiers and runs a plurality vote when predicting classes. These base classifiers are differentiated using a bootstrap technique to replicate the learning set [5]. The Random Forest (RF) and Pasting Small Votes ensemble models are a variation of bagging [6; 7]. The second group of ensemble models is boosting. Similar to bagging, boosting is also an aggregation of multiple classifiers. However, instead of using bootstrap, it uses the boosting technique to train the base classifiers. This technique can convert weak learning algorithms to achieve arbitrarily high accuracy [60]. Adaptive Boosting (AB) [24] and Gradient-Boosted Trees (GB) [25] are two well-known algorithms that use this technique. In stacked generalisation [67], the classifiers are stacked to one another so that the output of some first-level base classifiers becomes the input of a second-level meta classifier. The combination of expert methods [32] combines the weighted output of several base classifiers in the consecutive combiner.

While promising, we see the potential for improvement from these ensemble classifiers; they typically use only the same type of single classifier as the base classifier. To differentiate the output of each base classifier, they implement bootstrapping or boosting of the training sample inputs via bagging or boosting, stacking the classifiers, or applying different weights for their output. In contrast, our proposed MtCE combines several types of single classifiers working together as an ensemble model. As every single classifier has strengths and weaknesses in classification or detection, by combining different types of single classifiers in an ensemble, we use the strength of one classifier to 'cover' for the weakness of the other classifiers, and vice versa. The types of base classifiers ensembled in our MtCE can be customised by the user based on the problem and their experience and knowledge of this problem. To further increase the performance of our model, we implemented the bootstrap

technique [5] and the method to vote on output (majority or priority) to produce the final output. For the base classifiers of our proposed ensemble model, we implemented several standard single classifiers that performed well in previous research on text mining [9-12; 16; 17], including the Logistic Regression (LR) [46], Linear-kernel Support Vector Machine (LSVM) [13; 15], Multi-Layer Perceptron (MLP) [58], and Convolutional Neural Network (CNN) [40]. In addition, two other classifiers that are usually applied as the base classifiers in several ensemble models, namely Decision Tree (DT) [53] and Naïve Bayes (NB) [44], were also used as base classifiers. The MtCE was compared with ensemble models widely used in the literature, including Bagging Predictors (BP) [5], RF, AB and GB.

Theoretically, our work contributes to the artificial intelligence research domain, especially ML, by proposing a novel ensemble approach that can solve classification, detection and prediction problems. Most ensemble models in the literature either (1) have only one type of base classifier (such as RF, AB, BP and GB) or (2) have a small combination of fixed types and number of base classifiers that are designed to solve a particular problem [33; 55; 63; 64]. In contrast, our proposed model (the MtCE) provides the freedom to select the types and number of base classifiers depending on one's requirements. Bootstrapping is included in the MtCE process to improve the model's stability and accuracy. Finally, a priority threshold approach is introduced, in addition to majority voting, to decide the final result; this priority threshold approach can improve the recall and overcome the imbalanced issue.

From a practical perspective, this work also contributes to addressing online commerce security problems—we have successfully implemented a model that can detect fake online consumer reviews with better results than previous research. Product reviews from buyers are commonly used as a guide by most consumers to make purchasing decisions on electronic commerce these days [9]. Reliable and effective detection of fake reviews will increase the trustworthiness of online commerce [10]. On the other hand, sellers use consumer reviews to evaluate the brand perception and consumers' satisfaction [22]. For this reason, maintaining the quality of product reviews is vital. Therefore, fake review detection is an urgent task and a high priority for online commerce portal providers.

The rest of this paper is organised as follows. In the next section, we explain the design of our proposed model, the MtCE. After that, we discuss how we conducted experiments to test the performance of the MtCE, such as the datasets we used, testing framework, details of the textual-based approach, types of base classifiers, and the measurements. In Section 4, we discuss the results and analyses, and finally, we draw conclusions and outline the possibilities for future work.

## 2 RELATED WORK ON ENSEMBLE MODELS

The previous section provided a concise introduction to both fake review detection and ensemble models. This section discusses recent fake review detection studies from the literature (2016 to the present) that have proposed new ensemble models or implemented existing ensemble models for detection of fake reviews (see Table 1).

The objective of most previous studies has been to propose feature extraction approaches for the input of standard ML and deep learning (DL), including their ensemble models. As discussed in Section 1, these feature extraction approaches can either be textual-based, which creates features from the words/terms of the review text itself [11; 20; 21; 72]; or content-based, which creates features from the text and extracts features from the text's information, properties, sentiment polarities, and characteristics [2]; or behaviour-based, which emphasises the behaviour of the reviewers rather than their reviews [3; 19; 39; 45].

The above approaches have also been combined to improve the detection results [12; 26; 38; 61; 71]. These studies investigated existing ML, DL and ensemble models (i.e. AB, BP, RF, and GB) to detect fake reviews using features proposed via the combined approaches. The base classifiers of these ensemble models are of one type; for

example, the RF utilises decision trees as its base classifiers, and GB utilises regression trees—both of them utilise and modify the trees in their processes. Therefore, it is almost impossible to replace their tree-type base classifiers with others. While the type of the base classifier in some ensemble methods can be changed by design, the replacement must still be one type, e.g., it is impossible to mix more than one type of base classifier for AB and BP.

Some studies in the literature have also proposed novel ensemble models. For example, Sun et al. [63] proposed a bagging style of 2 SVMs and a CNN to detect fake reviews. They utilised special SVMs, namely BIGRAMS<sub>SVM</sub> and TRIGRAMS<sub>SVM</sub>, to detect terms features from the review text input, with a unique CNN model (Product Word Composition Classifier (PWCC)) to capture product-related review features. The output of these three classifiers were processed in a bagging style method to produce the final detection result. Similarly, a forest of decision trees, called the Neural Autoencoder Decision Forest, was proposed by Dong et al. [19] to detect fake reviews. Their ensemble was designed by combining the Deep Neural Decision Tree (proposed by Kotschieder et al. [37]) with the Autoencoder method. An integration of several CNNs and Gated Recurrent Neural Networks (GRNNs) were introduced by Ren and Ji [55] in 2017. The CNNs were implemented to produce continuous sentence vectors from words, then handled by several forward and backward GRNNs to produce document representations of the reviews. A unique model of CNN, namely bfGAN, was proposed by Tang et al. [64] in 2020, and combined with the SVM for fake review detection. They showed that their model can effectively detect cold-start spam/fake reviews; the cold start problem is when a new reviewer posts a review. Javed et al. [33] proposed an ensemble of three classifiers (CNNs) for fake review detection. Each classifier detects different feature groups: a textual-based, a content-based (non-textual), and a behaviour-based group of features; a voting mechanism is the used to produce the final detection. The above-discussed approaches are similar in terms of one important aspect: they have proposed a specific ensemble model with specific types and total number of base classifiers. In these ensemble models, the types and numbers of their base classifiers cannot be modified, since each base classifier has a specific task and purpose.

Our proposed ensemble model, the MtCE, is different. In the MtCE, diverse base classifiers can be mixed together so that the strength of one type of base classifier can overcome the weaknesses of other types of base classifier and vice versa. The total number of base classifiers and the number of each type of base classifier can also be configured freely based on specific requirements. Our model also applies the bootstrapping technique to ensure stability and improve accuracy. For the final classification result, we introduce a priority threshold technique that prioritises a particular class in conjunction with common majority voting that is usually applied by other ensembles. This priority threshold technique can overcome the imbalanced problem when applied to the scarce/rare class as well as improve the recall in binary classification if applied to the main (positive) class.

Table 1. An overview of related work (those algorithms in italic are ensemble models)

Author	Year	Featuring type	Algorithm
Sun et al. [63]	2016	Textual-based	<i>Bagging of 2 SVMs and PWCC</i>
Zhang et al. [71]	2016	Content- and behaviour-based	NB, SVM, DT, <i>RF</i>
Etaiwi and Naymat [21]	2017	Textual-based	NB, SVM, DT, <i>RF, GB</i>
Ren and Ji [55]	2017	Textual-based	SVM, GRNN, Recurrent Neural Network (RNN), CNN, <i>CNN-GRNN (Integrated)</i>
Dong et al. [19]	2018	Behaviour-based	<i>Neural Autoencoder Decision Forest</i>
Hazim et al. [26]	2018	Content and behaviour-based	<i>Generalised Boosted Regression Model (GBM) Gaussian, GBM Poisson, GBM Bernoulli, GB, AB</i>
Kumar et al. [39]	2018	Behaviour-based	LR, k-NN, NB, SVM, <i>RF, AB</i>
Zhang et al. [72]	2018	Textual-based	Recurrent CNN, SVM, CNN, <i>CNN-GRNN</i>
Barbado et al. [3]	2019	Behaviour-based	LR, DT, Gaussian NB (GNB), <i>RF, AB</i>

Martens & Maalej[45]	2019	Behaviour-based	DT, MLP, LSVM, RBF kernel SVM (RSVM), GNB, <i>RF</i>
Tang et al. [64]	2020	Behaviour-based	CNN, <i>CNN-bfGAN + SVM</i>
Alsubari et al. [2]	2020	Content-based	DT, <i>RF, AB</i>
Budhi et al. [11]	2021	Textual-based	LR, MLP, LSVM, <i>RF, AB, BP</i>
Budhi et al. [12]	2021	Content- and behaviour-based	DT, LR, LSVM, MLP, CNN, <i>RF, GB, AB, BP</i>
Elmogy et al. [20]	2021	Textual-based	LR, NB, k-NN, SVM, <i>RF</i>
Javed et al. [33]	2021	Textual-, content-, and behaviour-based	<i>Ensemble of 3 CNNs</i>
Shan et al. [61]	2021	Content- and behaviour-based	DT, SVM, NB, MLP, <i>RF</i>
Kumar et al. [38]	2022	Textual-, content-, and behaviour-based	Long short-term memory, Artificial Neural Network, RNN, RSVM, LR, k-NN, NB, <i>RF, GB, Light GB Method</i>

### 3 THE MTCE

The MtCE is a novel ensemble of classifiers developed in light of the fact that every single classifier has its strengths and weaknesses (see Figure 1). This ensemble is proposed based on the following idea:

*Every single classifier has been created with a different method, formulation and purpose, and thus, has strengths and weaknesses in different spots; for example, (1) for the same dataset, each classifier may correctly classify the different set of records; (2) a classifier is usually created to solve one or two problems, and not tested for the case of other problems; and (3) some classifiers perform well for smaller datasets but not for bigger ones, or vice versa. Therefore, by combining different single classifiers in one ensemble, the strengths of one classifier may ‘cover’ for the weaknesses of another.*

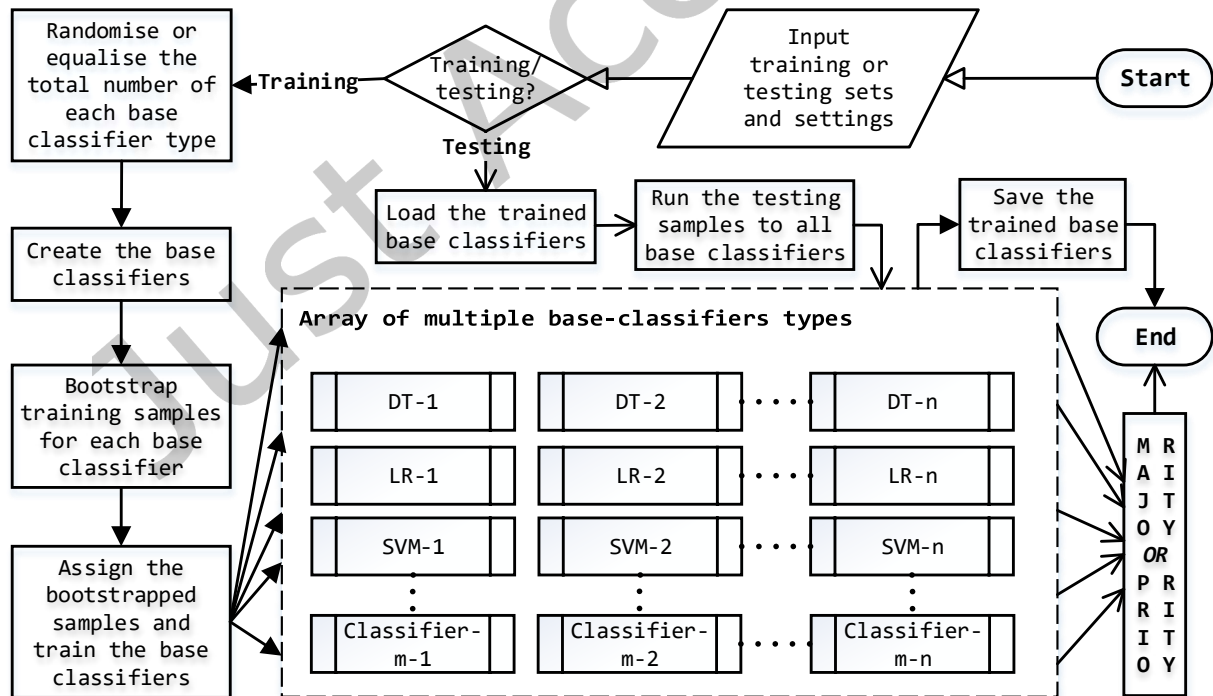


Figure 1: Design of the MtCE

In Figure 1, we depict two processes: the training process (the full head arrows  $\rightarrow$ ) and the testing process (the line head arrows  $\rightarrow$ ). All processes begin with inputting the training or testing sets and the parameter settings of the MtCE (see hollow head arrows  $\rightarrow$ ).

### 3.1 BASE CLASSIFIERS OF THE MTCE

Base classifiers here include several classifier objects that are utilised within an ensemble model. These classifiers will be run jointly in a specific way (depending on the ensemble design) to produce results. These results are then processed to produce the final result of the ensemble model. For our MtCE, first the user will decide on the types of base classifiers and their total numbers. After that, for the training process, the user determines if the base classifiers will be created in equally split numbers between the types setting or in random fashion. For example, suppose the total base classifiers are 10 of 3 types (LR, DT, NB); an equally split setting will create 4 LR, 3 DT and 3 NB, while for a random setting, each classifier type will be randomised from 1 to 8 (total classifiers – (total types – 1)) classifiers. While splitting the number of base classifier types equally is more sensible if the user does not know the exact nature of the problem at hand, randomising them sometimes could give better results. Moreover, with a simple modification in the implementation phase, the user can also set the exact different total number of each base classifier type, supposed they have a reason to do so. In this study, we did not test the effect of setting exact numbers of each type of base classifier since, for our problem, it will become similar to the randomised setting.

### 3.2 BOOTSTRAP SAMPLING

After creating the base classifiers, each base classifier is assigned a subset of training samples based on the bootstrap sampling process. Bootstrapping the samples could improve the stability and accuracy of the model [5]. The bootstrap sampling method draws sample data repeatedly with replacement from the source (data), based on the percentage setting [5; 35]. After the training process for each base classifier finishes, the weights and parameters of all base classifiers are saved in a file.

### 3.3 DETERMINING THE FINAL OUTPUT

The testing process is straightforward. After inputting the testing samples and loading the previously trained base classifiers, all samples are run with the base classifiers. The final output will be determined by one of two processes: the majority vote or the priority class threshold. The majority vote chooses the majority class outputs from the results of the base classifiers. If the majority votes of all classes are equal, the final result is the smallest class in the corpus. The priority class threshold provides a final result based on the priority class chosen in the setting; that is, if the positive class is chosen as the priority, then, if the positive class outputs from base classifiers are the same as or more than the threshold, the final output is a positive class. The priority class threshold can be set from absolute, which needs only one base classifier to pick the chosen class, to the maximum threshold before it becomes a majority vote (i.e. more than  $50\% + 1$  in a binary classification problem). This priority class threshold mechanism will focus on the class that is prioritised and will increase the performance of detection of this class.

## 4 TESTING THE PERFORMANCE OF THE MTCE

### 4.1 FRAMEWORK FOR TESTING

To evaluate our proposed ensemble, the MtCE, we implemented the comparison framework that we used previously to investigate classifiers for sentiment polarity detection of consumer reviews [10]. In this study, we modified the framework so that it could be applied to test the MtCE for fake review detection (see Figure 2). We used this framework to test different settings of the MtCE using similar datasets and comparing their performances. We ran the same test using several single classifiers as base classifiers of the MtCE and several well-known ensembles for comparison.

As can be seen in Figure 2, after a review text is loaded to a string, it is processed in the preprocessing subroutine to remove unnecessary words and corrections. Features are then extracted from the texts using a textual-based approach, as discussed in Section 1. After combining with the designated targets, these feature-target vectors are split into 10 folds for the cross-validation process. This same 10-fold setting is then trained and tested on several different combinations of MtCE models. Afterwards, the measurement results of these different MtCE models are compared to determine the best setting of the MtCE for fake review detection.

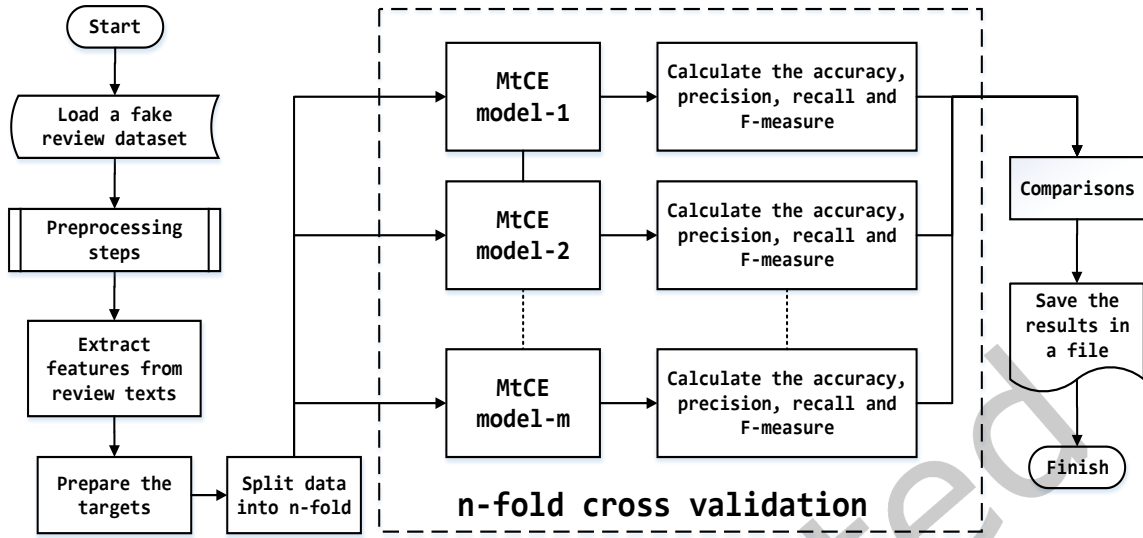


Figure 2: Framework to test and compare the MtCE

In the preprocessing steps (see Figure 3), we implemented several methods as follows:

- removal of punctuation, numbers and stopwords
- correction of spelling errors, elongation words and negative words
- POS tagging and lemmatisation.

Punctuation and number removal is a necessary step for a textual-based approach since the input features of this approach are the words. Therefore, we needed to remove the non-word parts of the text reviews. Stopword is a term for words commonly used in English sentences such as ‘the’, ‘is’, ‘at’, ‘which’ and ‘on’. These words can be removed from the text before it is used as a training record. These kinds of words may mislead detectors in recognising the trait words of a class because they are often the most common words in a corpus.

Misspelled words create unnecessary issues for detection since the detector will recognise them as different words from their correct forms. Like misspelled words, elongation words such as ‘yesss’, ‘fiine’ and ‘yoouu’ can also increase the diversity of words in the training example and make the classifiers harder to train. We utilised Peter Norvig’s code for spelling correction [48] to correct misspelled and elongation words. Negative words have many forms, depending on the grammar, but their purpose is similar—to change a positive sentence to be negative. Therefore, we shifted all negative words to their primary form to reduce the diversity of words.

POS tagging categorises words by their syntactic function, such as noun, pronoun, adjective, verb and adverb. This step is essential because it puts the words in their context [21] so that in the lemmatisation process, we can lemmatise the word to the correct context. Lemmatisation is a method for changing the word back to its basic form; POS lemmatisation returns the chosen word to its basic syntactic function. This step reduces the diversity of words inside the dataset and makes it easier to recognise.

Because we extracted the input features directly from social media messages, we used the BOW feature extraction method commonly used for textual-based features [18]. This method works by decomposing the entire text into a group of singular words. Similar to previous work [11], in this study, we used it in combination with the n-gram technique to capture singular words and terms that convey one meaning but are composed of multiple words. For terms, the BOW checks for the existence of a contiguous sequence of n words from the given text

sample. This study limited this checking to trigrams since sequences of more than three words are rare in real-world texts. After decomposing the text, the terms were sorted based on their frequency across the corpus.

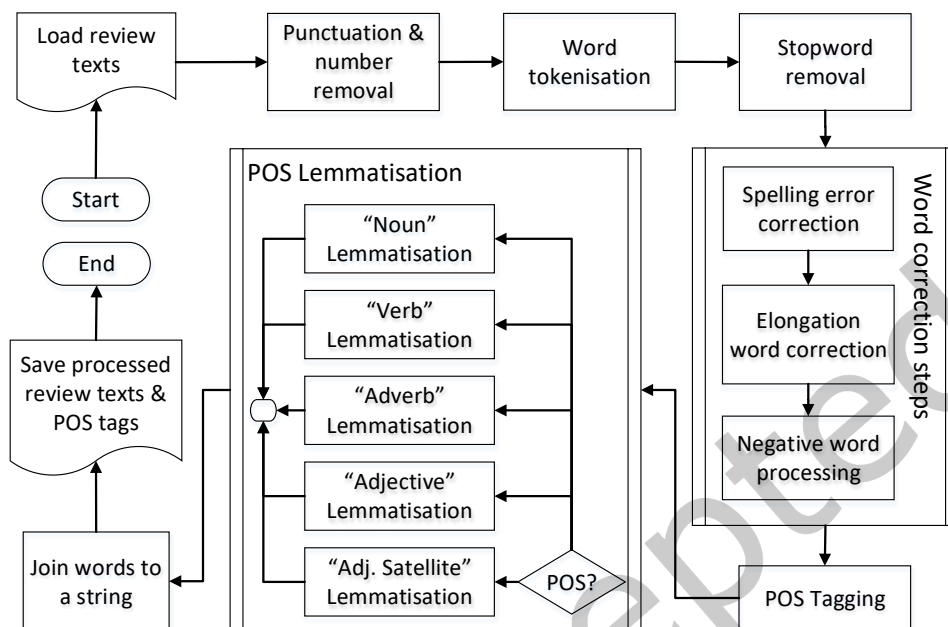


Figure 3: Preprocessing steps [11]

## 4.2 BASE CLASSIFIERS AND ENSEMBLE CLASSIFIERS FOR EXPERIMENTAL PURPOSE

As mentioned in Section 2, our proposed ensemble can apply multiple types of single classifiers. While in theory, any single classifier could be used as the base classifier of the MtCE, to test its performance, we investigated several combinations of ML and DL that performed well in previous research on text mining [9-12; 16; 17], which include the LR, LSVM, MLP, and CNN models. We also investigated two other classifiers, DT and NB, since these two models are commonly used as default base classifiers in some ensemble models, such as the RF [7], BP [5] and AB [24]. For comparison purposes, we tested some ensemble models BP, RF, AB and GB [25].

We built all ML classifiers and ensembles, except the CNN, using scikit-learn components [51]; the CNN was built with Keras [36] wrapped with scikit-learn components (scikit-learn does not provide a CNN component but a way to wrap DL components of Keras to be used with or within scikit-learn). To ensure the results could only be affected by our model implementation and not by modifying classifier parameters, we used the default parameters for all single classifiers used as base classifiers and the ensemble models.

### 4.3 DATASETS

Intuitively, our proposed ensemble classifier can be applied to a wide range of problems, similar to other ML ensemble classifiers. However, to test the performance of the MtCE, we conducted experiments using four public fake review datasets (see Table 2 for additional details). The public datasets from Ott et al. [49] and Li et al. [41] were created using domain experts, such as Amazon’s Mechanical Turk service (Turkers) and hotel employees, to provide false/fake reviews for some online services. Li et al.’s hotel dataset is an extension of Ott et al.’s dataset.

Table 2: Statistics for several public datasets

Dataset	Author	Domain	Total Records	Fake		Genuine	
				Total	%	Total	%



Ott et al. [49]	Hotel (Various)	1600	800	50.00	800	50.00
Li et al. [41]	Doctor (Various)	558	357	63.98	201	36.02
	Hotel (Various)	1880	1080	57.45	800	42.55
	Restaurant (Various)	402	202	50.25	200	49.75

#### 4.4 CROSS-VALIDATION AND MEASUREMENTS

We evaluated the performance of the MtCE using n-fold Cross-Validation (CV), which is widely applied to evaluate the generalisation performance of an algorithm [30], to reduce bias between the dataset and the training or testing set [54], and avoid overfitting [70]. We implemented scikit-learn [51] for performance measurements of our experiments. These measurements and their respective formulas can be found in Table 3.

Table 3: Measurement functions and formulas

No	Name	Sklearn Function	Equation
1	Accuracy	accuracy_score()	$A(y, \hat{y}) = \frac{1}{n_{samples}} \sum_{k=0}^{n_{samples}-1} 1(\hat{y}_k = y_k)$ <p>where <math>y</math> is the set of predicted pairs, <math>\hat{y}</math> is the set of true pairs and <math>n_{samples}</math> is total samples.</p>
2	Precision	precision_score()	$P(y_k, \hat{y}_k) = \frac{tp}{tp + fp}$ <p>where <math>k</math> is the set of classes, <math>y_k</math> is the subset of <math>y</math> with class <math>k</math>, <math>tp</math> is true positive, and <math>fp</math> is false positive.</p>
3	Recall	recall_score()	$R(y_k, \hat{y}_k) = \frac{tp}{tp + fn}$ <p>where <math>fn</math> is false negative.</p>
4	F-measure/F1	f1_score()	$F_1(y_k, \hat{y}_k) = 2 * \frac{P(y_k, \hat{y}_k) * R(y_k, \hat{y}_k)}{P(y_k, \hat{y}_k) + R(y_k, \hat{y}_k)}$

## 5 RESULTS AND DISCUSSIONS

### 5.1 EXPERIMENTING WITH THE MTCE FOR FAKE REVIEW DETECTION

The first group of experiments aimed to test the performance of the MtCE for fake review detection. For its base classifiers, we implemented the CNN, LR, LSVM and MLP, which were performing well in our previous studies [9-12]. We also implemented the DT and NB, which were used as default base classifiers in many ensembles (e.g. RF, BP, and AB). The parameters of these base classifiers are the defaults of the scikit-learn components used to implement them [51]. As shown in the Appendix, we tested all possibilities from 2-, 3-, 4- to 5-combination and presented 11 selected combinations in Table 4. Besides combining base classifiers, the other settings were fixed: total base classifiers = 15, shared equally for each type; bootstrap = 0.5; and final output = majority vote. The datasets we used in the experiments were Ott et al.'s dataset [49] and all three of Li et al.'s datasets [41]. All experiments were conducted using the 10-fold CV method. For the comparison, we also tested the datasets using single classifiers as above and several well-known ensemble classifiers (RF, BP, AB and GB). Results of the single and ensemble classifiers are provided in Table 5.

Table 4: Results of the combination of base classifiers for the MtCE model (selected)

Num <sup>1</sup>	MtCE Base Classifier Combination	Ott et al.'s Dataset <sup>2</sup>				Li et al.'s Dataset (Doctor) <sup>2</sup>			
		Acc	Pre	Rec	F1	Acc	Pre	Rec	F1
2	LR-MLP	<b>88.00</b>	<i>86.85</i>	<b>89.44</b>	<b>88.08</b>	84.41	<i>84.23</i>	<b>92.97</b>	88.33
5	LSVM-MLP	87.50	85.72	<b>90.02</b>	<b>87.74</b>	<b>86.56</b>	<b>87.05</b>	<b>93.12</b>	<b>89.88</b>
7	LSVM-NB	<b>89.19</b>	<b>88.32</b>	<b>90.43</b>	<b>89.31</b>	<i>84.75</i>	86.85	<b>90.32</b>	88.19
9	MLP-NB	<b>89.63</b>	<b>88.52</b>	<b>91.22</b>	<b>89.80</b>	<i>85.31</i>	<i>84.00</i>	<b>95.30</b>	<i>89.12</i>
18	LR-LSVM-NB	<b>89.00</b>	<b>88.68</b>	<b>89.59</b>	<b>89.07</b>	84.58	85.70	<b>91.44</b>	88.34
22	LSVM-MLP-NB	<b>89.13</b>	87.86	<b>91.00</b>	<b>89.31</b>	85.13	86.39	91.88	88.79
26	CNN-LR-MLP	<b>88.19</b>	<b>87.32</b>	<b>89.33</b>	<b>88.27</b>	84.40	83.75	<b>93.51</b>	88.29
35	LR-LSVM-MLP-NB	<b>89.19</b>	<b>88.77</b>	<b>89.95</b>	<b>89.28</b>	85.66	87.01	91.38	88.94
37	LR-MLP-DT-NB	<b>88.75</b>	<b>88.41</b>	89.25	<b>88.79</b>	84.06	83.85	<b>93.06</b>	88.04
47	LR-LSVM-MLP-DT-NB	<b>88.94</b>	<b>88.38</b>	<b>89.56</b>	<b>88.94</b>	84.24	84.89	92.00	88.12
50	CNN-LR-LSVM-MLP-NB	<b>88.88</b>	<b>88.21</b>	<b>89.83</b>	<b>88.95</b>	86.21	86.45	<b>93.18</b>	89.54
Li et al.'s Dataset (Hotel) <sup>2</sup>					Li et al.'s Dataset (Restaurant) <sup>2</sup>				
2	LR-MLP	<b>87.34</b>	85.09	<b>94.30</b>	<b>89.45</b>	85.11	83.99	<b>88.36</b>	85.69
5	LSVM-MLP	86.12	84.71	<b>92.67</b>	<b>88.44</b>	<b>85.83</b>	83.98	<b>88.55</b>	<b>86.02</b>
7	LSVM-NB	<b>87.45</b>	87.09	<b>91.75</b>	<b>89.34</b>	85.35	84.51	86.33	85.21
9	MLP-NB	87.23	86.69	<b>91.92</b>	<b>89.17</b>	86.04	85.36	87.18	<b>86.15</b>
18	LR-LSVM-NB	85.96	85.20	<b>91.34</b>	88.14	<b>86.58</b>	85.71	<b>88.08</b>	<b>86.41</b>
22	LSVM-MLP-NB	86.97	86.09	<b>92.34</b>	<b>89.05</b>	85.32	83.68	<b>88.65</b>	85.75
26	CNN-LR-MLP	85.85	84.79	<b>91.96</b>	88.17	83.34	82.42	84.09	82.88
35	LR-LSVM-MLP-NB	86.70	85.12	<b>92.92</b>	88.78	82.32	81.51	84.77	82.80
37	LR-MLP-DT-NB	<b>87.82</b>	86.67	<b>93.26</b>	<b>89.81</b>	84.38	83.61	85.67	84.57
47	LR-LSVM-MLP-DT-NB	<b>87.55</b>	86.19	<b>93.28</b>	<b>89.58</b>	84.63	83.29	87.36	85.06
50	CNN-LR-LSVM-MLP-NB	86.76	85.98	<b>91.98</b>	<b>88.81</b>	83.80	82.59	86.17	83.98

<sup>1</sup> The number here is associated with the number in Appendix

<sup>2</sup> **Bold-green** = the MtCE result is higher than that of all base classifiers in Table 5; *Italic-red* = the MtCE result is higher than all base classifiers; Normal-black = at least one base classifier result is higher than that of the MtCE.

Table 5: Results of single and ensemble classifiers

Num.	Classifiers	Ott et al.'s Dataset				Li et al.'s Dataset (Doctor)			
		Acc	Pre	Rec	F1	Acc	Pre	Rec	F1
1	CNN	79.63	81.48	77.36	79.13	67.76	70.64	88.40	77.02
2	LR	87.13	87.05	87.25	87.11	83.35	85.48	89.71	87.26
3	LSVM	85.88	85.67	86.21	85.89	83.13	86.28	87.80	86.91
4	MLP	87.56	87.05	88.37	87.66	85.67	85.89	92.94	89.16
5	DT	71.50	72.21	69.71	70.83	65.79	72.98	74.01	73.32
6	NB	88.50	87.97	89.45	88.63	87.12	90.68	88.98	89.68
7	RF	87.00	87.44	86.62	86.92	76.35	74.05	97.24	83.98
8	BP	75.00	74.23	76.79	75.41	74.19	74.50	90.21	81.45
9	AB	80.06	79.64	80.82	80.09	74.93	79.56	81.28	80.16

10	GB	82.81	83.67	81.63	82.56	76.52	77.08	90.63	82.99
		Li et al.'s Dataset (Hotel)				Li et al.'s Dataset (Restaurant)			
1	CNN	78.83	82.48	80.63	81.29	71.63	71.66	73.20	71.75
2	LR	85.59	85.50	90.25	87.77	81.59	81.03	84.52	81.74
3	LSVM	84.20	84.59	88.26	86.33	82.11	79.93	84.40	81.53
4	MLP	86.12	85.89	90.82	88.24	85.56	84.85	87.69	85.73
5	DT	68.88	73.81	71.15	72.38	60.24	60.34	63.27	61.11
6	NB	87.29	89.73	87.92	88.78	86.08	86.36	86.63	86.10
7	RF	81.97	80.78	90.28	85.17	81.34	78.31	87.86	81.88
8	BP	75.96	76.35	84.14	80.03	71.65	71.08	74.51	71.95
9	AB	79.89	80.77	84.87	82.67	70.89	70.64	70.35	70.38
10	GB	80.48	78.83	89.94	84.00	76.12	77.43	75.46	75.78

Comparing Table 4 with Table 5, some combinations of base classifiers in the MtCE exceed the performance of all of the base classifiers (see the bold-green scores in Table 4). For example, accuracy of the MtCE(LR-MLP) for Ott et al.'s = 88%; in comparison, in Table 5, the accuracy values of LR and MLP are 87.13% and 87.56%. However, not all combinations improved and supported each other as we hoped. A few weak classifiers degraded the performance of stronger classifiers when combined with these. This gave rise to the result that the performance of the MtCE was in the middle of the performance of its base classifiers; for example, the LSVM accuracy for Ott et al.'s = 85.88%, while MLP = 87.56%, and the combination of both in MtCE(LSVM-MLP) = 87.50%. This implies that the LSVM degraded the performance of MLP, and given this, rather than implementing our proposed ensemble, it would be better to use the MLP directly. However, in the same case, MtCE(LSVM-MLP) for Ott et al.'s improved the recall, from 86.21% (LSVM) and 88.37% (MLP) to 90.02%. In a case such as this, using our ensemble would be acceptable since the recall in binary classification/detection is the same as the accuracy of the positive class or the main class (in our problem, the fake review class). Thus, high recall means the ability of the ensemble to detect positive classes is also high.

Another finding from this set of experiments is that different datasets might need different base classifier combinations. For example, MtCE(MLP-NB) that performed best for Ott et al.'s dataset performed worst for Li et al.'s doctor dataset (increasing the recall but decreasing all other measurements). The best combination for Li et al.'s doctor dataset was MtCE(LSVM-MLP); for Li et al.'s hotel dataset was MtCE(LR-MLP-DT-NB); and for Li et al.'s restaurant dataset was MtCE(LR-LSVM-NB). A closer inspection of Table 5 shows that the base classifiers of MtCE's best combination for a particular dataset mostly performed well on the same dataset too, when they were used in stand-alone modes. For example, on Ott's MtCE(MLP-NB), the MLP and NB also performed well on Ott's dataset, although they were not as good as when combined in the MTCE. This observation is valid on other combinations for other datasets, except the DT on Li et al.'s hotel dataset, which was not performing well by itself but performed best in combination with other base classifiers in the MtCE. This indicates that, while combining good classifiers could produce better results, infusing a weak classifier sometimes could also boost the results of MtCE.

We found that the MtCE performed better than other ensembles of homogenous classifiers, as indicated by comparing the results in Table 5 (numbers 7–10) to the MtCE results in Table 4. These facts suggest that, when performed correctly, the combination of multi-type classifiers could cover the weakness of each and outperform the combination of a homogenous type of ensemble classifiers such as the RF, BP, AB or GB. However, for the cases of RF, BP and AB, we suspect the cause is also because these ensembles use the DT as their base

classifiers/predictors. DT performance was not good on all datasets, which therefore affected the performance of its ensemble variants.

While the deep learning-based CNN as a stand-alone classifier did not perform as well as other single classifiers, somewhat surprisingly, it performed exceptionally well when combined with other classifiers in the MtCE. This increase in accuracy was significant; for example, MtCE(CNN-LR-LSVM-MLP-NB) is 7.93% for Li et al.'s hotel dataset to 18.45% for Li et al.'s doctor dataset. The improvement in recall was also good, with a minimum of 4.78% for Li et al.'s doctor dataset to 12.97% for Li et al.'s restaurant dataset. Therefore, the deep learning-based CNN could be combined perfectly with other single ML classifiers as base classifiers for the MtCE; it gave good results and improved all other base classifiers' performances as well. This is another strong indication that combining different classifiers with their strengths and weaknesses in our proposed MtCE can improve the overall performance. It is because the weaknesses of one base classifier can well be covered by the strengths of other base classifiers (i.e. ensemble diversity).

## 5.2 EFFECTS OF MTCE PARAMETERS

In this section, we discuss several input parameters of the MtCE that may affect the performance of this proposed model. All of the settings are the same as in Section 5.1 except the parameter we are testing. For these experiments, we chose four base classifier combinations of the MtCE in Table 4—that is, MtCE(LSVM-MLP), MtCE(LSVM-NB), MtCE(LR-LSVM-NB), MtCE(LR-MLP-DT-NB) and MtCE(LR-LSVM-MLP-DT-NB)—as these are the five combinations of base classifier types that have performed the best across all four datasets. We ran all experiments for the parameter testing using Ott et al.'s dataset only, since these five MtCE combination settings perform well in this dataset (see the bold-green scores in Table 4). However, after finding the ideal set of parameters, we ran them with the other datasets and compared the results with our previous approach implemented on the same datasets.

### 5.2.1 TOTAL NUMBER OF BASE CLASSIFIERS

First, we investigated the effect of base classifiers' total numbers. Here, we ran experiments on all five combinations with the total number of base classifiers varying from 5 to 100. Results can be seen in Figure 4 for accuracy, Figure 5 for precision, Figure 6 for recall and Figure 7 for F-measure. From these figures, we can see that accuracy, precision, recall and F-measure were increasing at first, then, after 20 classifiers, all scores fluctuated around a number  $\pm 1.5\%$ .

Accuracy fluctuated at 88.5%, precision at 88% and recall at 90%, except for MtCE(LSVM-MLP), which was 1% lower.

Intuitively, these fluctuations after 20 classifiers mean that adding more classifiers would not improve the measurements a lot. However, since the case is only for fake review detection on Ott et al.'s dataset, we will let the user decide on the number of base classifiers needed for their problem. Next, since we think a 1.5% difference is quite significant, we used the best accuracy results for each MtCE combination for the next batch of experiments. These combinations are 20 classifiers for MtCE(LSVM-MLP), 55 classifiers for MtCE(LR-LSVM-NB), 80 classifiers for MtCE(LR-MLP-DT-NB) and 85 classifiers for MtCE(LR-LSVM-MLP-DT-NB).

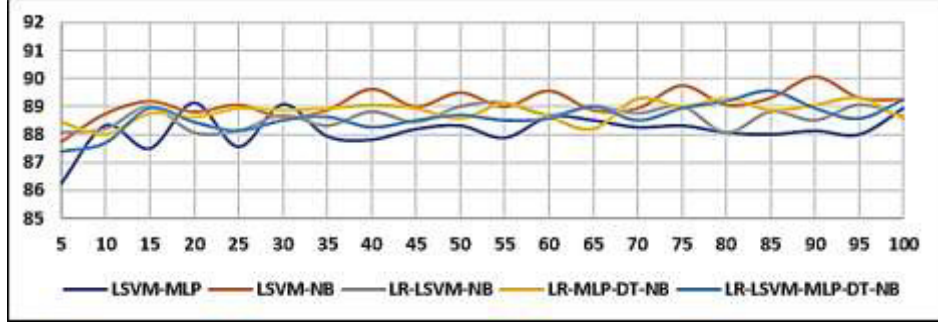


Figure 4: The accuracy (y-axis, in %) of the MtCE with 5–100 total base classifiers (x-axis)

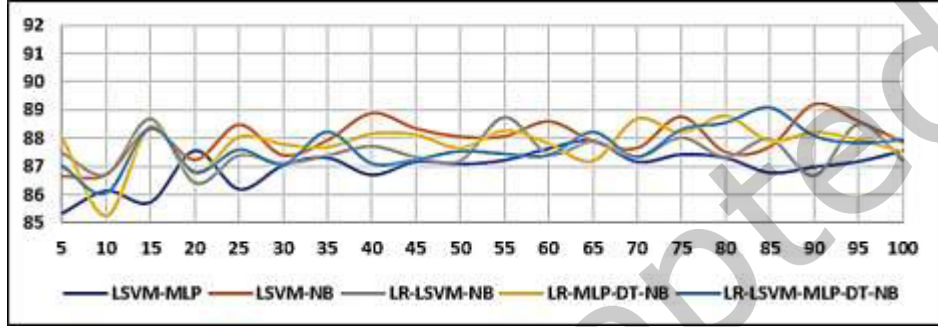


Figure 5: The precision (y-axis, in %) of the MtCE with 5–100 total base classifiers (x-axis)

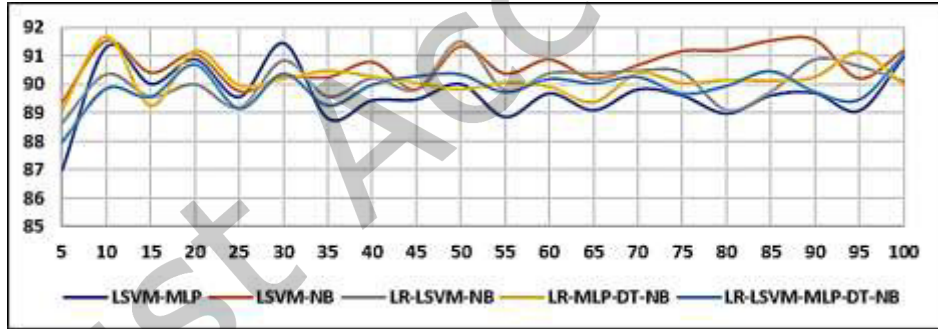


Figure 6: The recall (y-axis, in %) of the MtCE with 5–100 total base classifiers (x-axis)

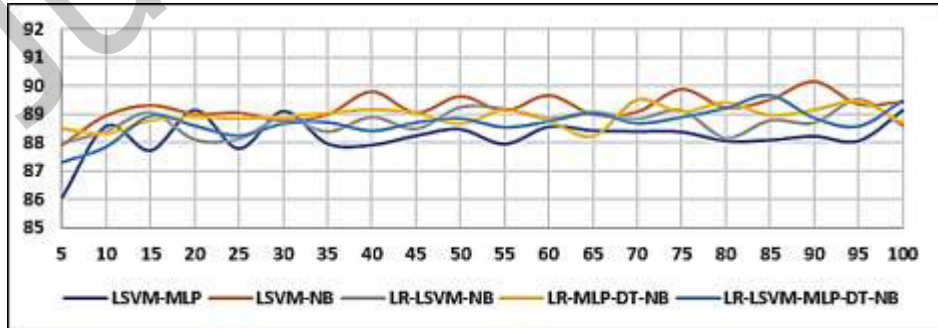


Figure 7: The F-measure (y-axis, in %) of MtCE with 5–100 total base classifiers (x-axis)

### 5.2.2 EQUALLY SPLIT OR RANDOMLY ASSIGNED BASE CLASSIFIER TYPE

To investigate the equally split vs randomly assigned base classifier type, we ran the experiments 10 times for each equally split and randomly assigned setting for each MtCE combination. The average results and their standard deviation can be seen

in Table 6. This table shows that the equally split option gave slightly better average results for all measurements relative to the random option. This implies that when the total of each base classifier is equal, they support each other, and the strength of one type of classifier covers the weakness of the other type when combined in the MtCE. On the other hand, the random option could make one or two base classifiers more numerous than the others, so that they dominated the detection process. The standard deviation of all measurement scores below 1 means the variation of the results from 10 runs was low and close to each other. In more detail, we can see that the standard deviation of 2-combination, the MtCE(LSVM-MLP), on the equally split setting was lower than the random assignment. This is expected since, in an equally split setting, each base classifier total is the same for all 10 runs, while in the random option, this varies. However, the exact opposite happened in 5-combination, the MtCE(LR-LSVM-MLP-DT-NB). Here, the standard deviation of the equally split setting was higher than the random option, implying that the randomness might make the base classifiers more homogenous than split equally between its 5 base classifier types. For the 3- and 4-combinations, the standard deviation scores are similar.

Table 6. Results of equally split vs randomly assigned type of MtCE base classifiers on 10 runs of the 10-fold CV process.

Base Classifier Combination	Equally Split or Randomly Assigned	Accuracy		Precision		Recall		F-measure	
		Avg.	St.Dev	Avg.	St.Dev	Avg.	St.Dev	Avg.	St.Dev
LSVM-MLP	Equal	88.21	0.38	86.76	0.35	90.23	0.56	88.39	0.41
	Random	87.82	0.60	86.41	0.59	89.81	0.69	88.00	0.60
LSVM- NB	Equal	89.44	0.36	88.20	0.50	91.12	0.39	89.57	0.37
	Random	89.37	0.19	88.16	0.38	90.93	0.35	89.47	0.22
LR-LSVM-NB	Equal	88.97	0.35	88.14	0.48	90.06	0.42	89.04	0.35
	Random	88.72	0.37	87.79	0.41	89.89	0.49	88.78	0.38
LR-MLP-DT-NB	Equal	89.18	0.32	88.34	0.44	90.24	0.39	89.21	0.29
	Random	89.03	0.35	88.19	0.46	90.12	0.38	89.09	0.35
LR-LSVM-MLP-DT-NB	Equal	88.28	0.66	87.37	0.64	89.59	0.86	88.41	0.67
	Random	88.12	0.37	87.10	0.39	89.49	0.49	88.21	0.39

### 5.2.3 BOOTSTRAP EFFECT

The subsequent experiments aimed to investigate the bootstrap parameter. We used fixed parameters as above, except the bootstrap setting ranged from 0.25 to 1 (no bootstrapping). As can be seen in Figure 8, the best accuracy, precision and recall were often achieved when the bootstrap setting was 0.5, where every base classifier was trained using 50% of training samples randomly. This indicates that bootstrap setting 0.5 is ideal for the tested MtCE settings, i.e., the five best base-classifier combinations. Therefore, intuitively, it can be used as the default bootstrap setting of MtCE. Based on these findings, for the next set of experiments, we set the bootstrap to be 0.5.

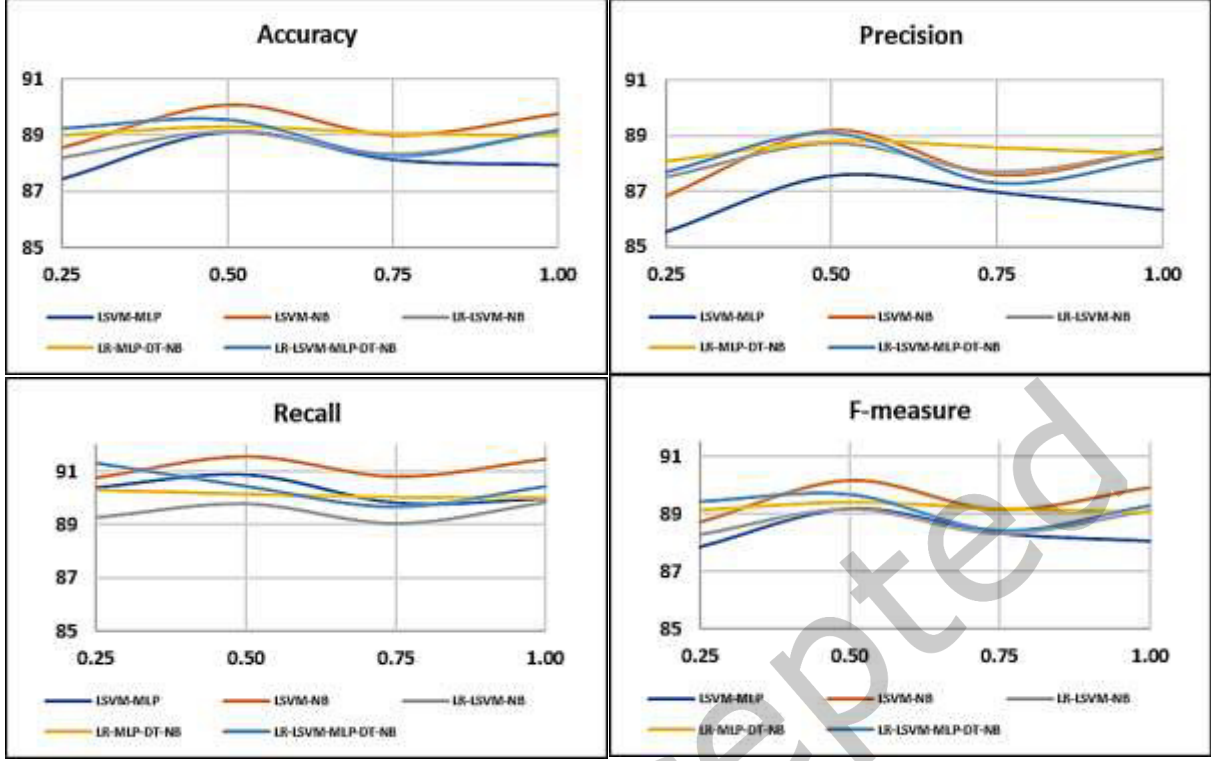


Figure 8: The accuracy, precision, recall and F-measure (y-axis, in %) of the MtCE with bootstrap settings from 0.25 to 1.00 (x-axis)

#### 5.2.4 FINAL OUTPUT: MAJORITY VOTE OR CLASS PRIORITY THRESHOLD

The last parameter to test is how the output of base classifiers should be processed. There are two options in terms of transferring the output of base classifiers to the final output: majority vote and priority class. To test the effect of priority, we set the priority to the positive class (fake class), in the range between absolute priority and 45% priority. As mentioned above, absolute priority means the final result will be recorded as positive/fake if one or more of the base classifiers records this. Percentage priority means the final result will be recorded as positive/fake if the number of base classifiers with this result is the same or more than the percentage threshold. The results of the experiments can be seen in Figure 9.

As per Figure 9, the priority setting increases the recall but decreases other measurements. Moreover, as noted, the higher the recall, the more accurate the positive class (fake class) detection. Thus, we can conclude that the priority setting is helpful when samples of positive class are scarce or when the focus of research is on anomaly detection. The key is how high we choose the percentage of priority so that the increase in recall does not greatly decrease the other measurements. For example, in Figure 9, for 25% priority of MtCE(LSVM-MLP), where recall increased by 5.05%, the accuracy, precision and F-measure decreased by 2.75%, 7.15% and 1.70%, respectively.

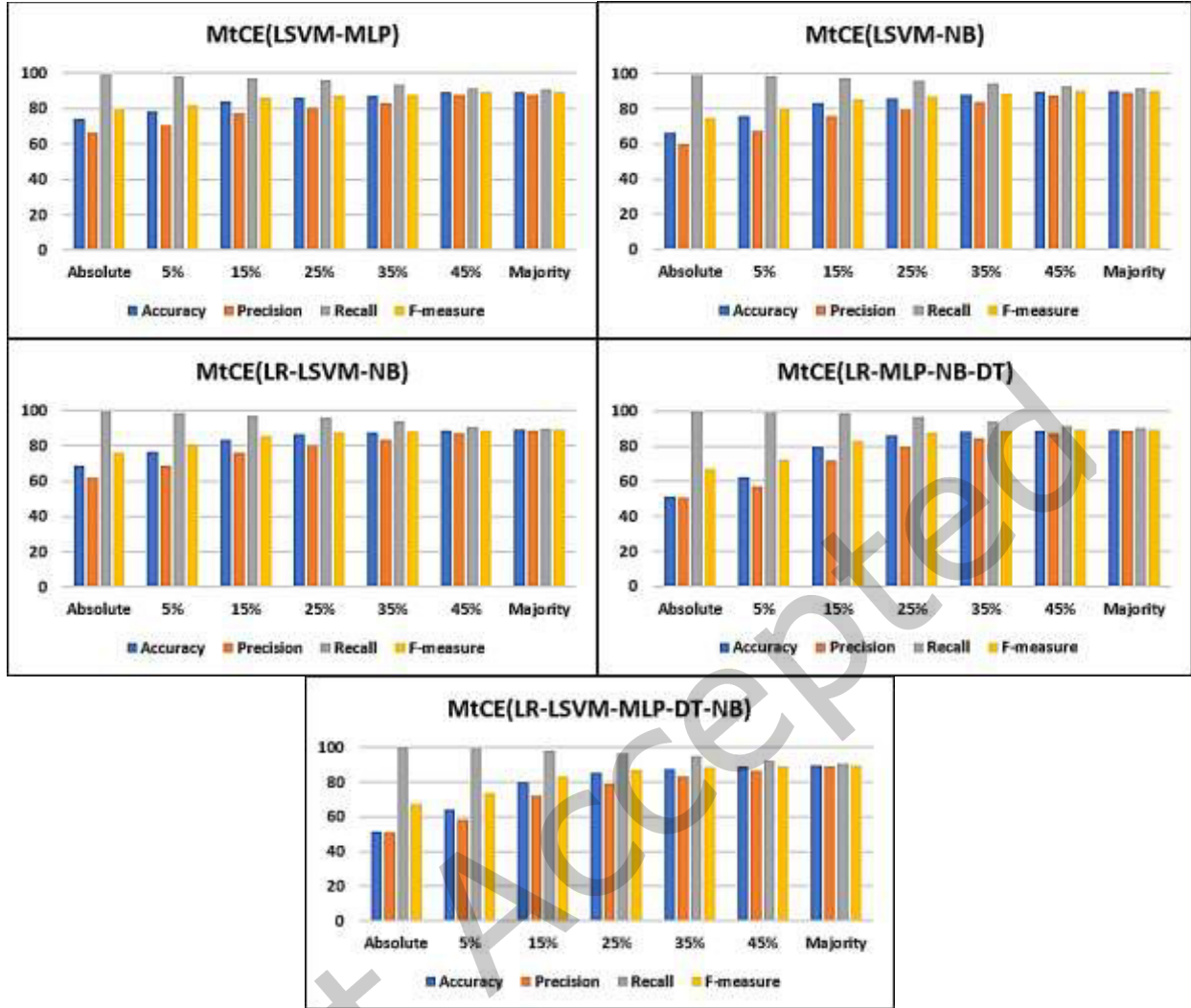


Figure 9: Results of final target experiments, in terms of accuracy, precision, recall and F-measure (y-axis, in %), from absolute priority to 45% priority, and majority (x-axis)

### 5.3 COMPARISON TO OTHER STUDIES ON THE SAME DATASETS

As commonly presented in related studies [42; 55; 57; 72], in this section, we present the MtCE's and other models' best results in light of various considerations. Many factors can influence results, such as the measurement formulas or tools, datasets used, number of records involved in experiments, and how the experiments were conducted (e.g., CV vs traditional train-test-validation split). Therefore, the results presented in Table 7 should be read with the following considerations in mind:

1. We present the MtCE's and other studies' results based on the same datasets (see Table 2). We used all samples that could be processed from each dataset. Note that some studies used only a subset of the dataset or split the dataset into subsets and measured each subset differently. All of our experiments were conducted using 10-fold CV.
2. It is impossible to ensure that all the studies used the same formulas to measure accuracy, precision, recall and F-measure. Therefore, MtCE results were measured using widely used formulas for binary target prediction (see Table 3). All measurements were in the 'macro' average setting using scikit-learn measurement components [51].
3. The original results from the dataset creators were used as the baseline (see the underlined description and scores in Table 7. We present only the best result for each dataset from each study.



Table 7: Best performance of the MtCE and other studies for the same datasets

No	Dataset	Description	Acc (%)	Pre (%)	Rec (%)	F1 (%)
1	Hotel (various) [49]	<u>Ott et al. [49], SVM on positive sentiment subset</u>	<u>88.4</u>	<u>89.1</u>	<u>87.5</u>	<u>88.3</u>
		Fusilier et al. [27], PU-L modified	-	85.2	72.8	78
		Rout et al. [57], DT	92.11	-	-	-
		Etaiwi and Naymat [21], SVM, all preproc. steps	85	51	86	-
		Zhang et al. [72], DRI-RCNN, 0.9 T.P.	-	-	-	86.59
		Budhi et.al. [12], GB, over-sampling	70.62	67.58	81.29	73.8
		<i>MtCE(LSVM-NB, 90, Equal, 0.5, Majority)*</i>	<i>90.06</i>	<b><i>89.19</i></b>	<b><i>91.56</i></b>	<b><i>90.16</i></b>
2	Hotel (various) [41]	<u>Li et al. [41], OvR SAGE-Unigram</u>	<u>81.8</u>	<u>81.2</u>	<u>84</u>	-
		Ren and Ji [55], Integrated, Employee/Turker subset	92.6	-	-	90.1
		Li et al. [42], SWNN	-	84.1	87	85
		Budhi et.al. [12], GB, under-sampling	71.29	73.61	82.79	77.93
		<i>MtCE(LR-MLP-DT-NB, 15, Equal, 0.5, Majority)*</i>	<i>87.82</i>	<b><i>86.67</i></b>	<b><i>93.26</i></b>	<i>89.81</i>
		<u>Li et al. [41], OvR SAGE-Unigram</u>	<u>81.7</u>	<u>84.2</u>	<u>81.6</u>	-
3	Restaurant (various) [41]	Ren and Ji [55], Integrated, Customer/Turker subset	86.9	-	-	86.8
		Li et al. [42], SWNN	-	83.3	88.2	81
		Budhi et.al. [12], GB, over-sampling	72.44	71.23	75.16	73.14
		<i>MtCE(LR-MLP-DT-NB, 80, Equal, 0.5, Majority)*</i>	<i>86.07</i>	<b><i>84.47</i></b>	<b><i>88.89</i></b>	<i>86.37</i>
		<u>Li et al. [41], OvR SAGE-Unigram</u>	<u>74.5</u>	<u>77.2</u>	<u>70.1</u>	-
		Ren and Ji [55], Integrated, Customer/Turker subset	76	-	-	74.1
4	Doctor (various) [41]	Li et al. [42], SWNN	-	83.7	87.6	82.9
		Budhi et.al. [12], GB, over-sampling	73.35	79.44	86.94	83.02
		<i>MtCE(LSVM-MLP, 15, Equal, 0.5, Majority)*</i>	<b><i>86.56</i></b>	<b><i>87.05</i></b>	<b><i>93.12</i></b>	<b><i>89.88</i></b>

\*MtCE(<type of base classifiers>, <total base classifiers>, <equally split or random assign of base classifiers>, <bootstrap percentage>, <majority or priority output>)

We do not claim that the MtCE is better or worse than methods in other studies since several factors can affect results. However, we can confidently compare the MtCE to our previous study in fake review detection [12] because the measurement formulas are similar. As is evident in Table 7, compared with our previous approach, the MtCE is superior. Accuracy improves from a minimum of 13.2% in Li et al.'s doctor dataset to a maximum of 19.4% in Ott et al.'s dataset; precision improves from 7.6% to 21.6% and recall from 6.2% to 13.7%.

## 6 CONCLUSION AND FUTURE WORK

From the experiments above, we can conclude that our proposed ensemble, the MtCE, was able to adequately detect fake or fraudulent reviews, which have become an acute problem on e-commerce platforms. Compared with our previous approach, the MtCE improved accuracy up to 19.4%, precision up to 21.6% and recall up to 13.7%.

Not all combinations of base classifier types produced the expected result of an improvement in the performance on all base classifiers of the MtCE. In some combinations, weak classifiers dragged down the performance of stronger classifiers, especially in terms of accuracy measurements. However, when the combinations were correct, the MtCE produced better results than its base classifiers in stand-alone modes. It is important to note that while accuracy was not on top, recall was the highest in many cases. This is a positive result since, in binary classification, recall is the same as accuracy of a positive case. In this case, the MtCE was able to better detect the fake review class. Comparison with well-known ensembles, such as the RF, AB, BP and GB, showed that some correct combination of the MtCE could outperform the other ensembles, proving that combining multi-type classifiers in one ensemble process performed better than combining the same classifiers, as is usually done in other ensemble methods.

We proved that DL methods such as the CNN model could be combined with ML counterparts as base classifiers to give better results than if run alone. The number of base classifiers affected performance detection; however, accuracy and other measurements oscillated after 20 base classifiers. While reducing the amount of data for the training process, the bootstrap setting could also affect the performance of the MtCE. From the experiments, we found that 50% bootstrapping gives better results than other percentages, including the non-bootstrap (bootstrap setting 100%). While the majority vote setting for the output offers a fair chance for each class to be detected, the priority threshold boosts detection of the prioritised class. This would be useful if the dataset is imbalanced or the MtCE is used to detect/classify anomalies.

Despite the extensive experimental studies presented in this paper, there remain opportunities/possibilities for further improvement. First, we are aware that multiple types of base classifiers will increase the complexity of the ensemble and increase the cost and time processing for parameter tuning. To overcome this problem, in the near future, we plan to expand on our previously proposed algorithms for ensemble parameter optimisation, namely the Multi-level Particle Swarm Optimisation (PSO) and its parallel version [8]. These nature-inspired algorithms are based on a low-cost PSO algorithm. They were designed to optimise the parameters of an ensemble model, such as BP and AB, choose its base classifier type, and optimise the parameters of these base classifier candidates. Other avenues for future work include investigating the implementation of the MtCE for multi-class problems, heavily imbalanced datasets, and anomaly detection.

## REFERENCES

- [1] AKRAM, A.U., KHAN, H.U., IQBAL, S., IQBAL, T., MUNIR, E.U., and SHAFI, M., 2018. Finding rotten eggs: A review spam detection model using diverse feature sets. *KSII Transactions on Internet and Information Systems* 12, 10. DOI= <http://dx.doi.org/10.3837/tis.2018.10.026>.
- [2] ALSUBARI, S.N., SHELKE, M.B., and DESHMUKH, S.N., 2020. Fake reviews identification based on deep computational linguistic features. *International Journal of Advanced Science and Technology* 29, 8s, 3846-3856.
- [3] BARBADO, R., ARAQUE, O., and IGLESIAS, C.A., 2019. A framework for fake review detection in online consumer electronics retailers. *Information Processing & Management* 56, 4, 1234-1244. DOI= <http://dx.doi.org/10.1016/j.ipm.2019.03.002>.
- [4] BING, L., WONG, T.-L., and LAM, W., 2016. Unsupervised extraction of popular product attributes from e-commerce web sites by considering customer reviews. *ACM Transactions on Internet Technology* 16, 2, 1-17. DOI= <http://dx.doi.org/10.1145/2857054>.
- [5] BREIMAN, L., 1996. Bagging predictors. *Machine Learning* 24, 2, 123-140. DOI= <http://dx.doi.org/https://doi.org/10.1007/bf00058655>.
- [6] BREIMAN, L., 1999. Pasting Small Votes for Classification in Large Databases and On-Line. *Machine Learning* 36, 1 (1999/07/01), 85-103. DOI= <http://dx.doi.org/10.1023/A:1007563306331>.
- [7] BREIMAN, L., 2001. Random forests. *Machine Learning* 45, 1, 5-32. DOI= <http://dx.doi.org/https://doi.org/10.1023/a:1010933404324>.
- [8] BUDHI, G.S., CHIONG, R., and DHAKAL, S., 2020. Multi-level particle swarm optimisation and its parallel version for parameter optimisation of ensemble models: a case of sentiment polarity prediction. *Cluster Computing* 23, 3371-3386. DOI= <http://dx.doi.org/https://doi.org/10.1007/s10586-020-03093-3>.
- [9] BUDHI, G.S., CHIONG, R., PRANATA, I., and HU, Z., 2017. Predicting rating polarity through automatic classification of review texts. In *Proceedings of the 2017 IEEE Conference on Big Data and Analytics (ICBDA)*, Kuching, Malaysia, 19-24. DOI= <http://dx.doi.org/https://doi.org/10.1109/ICBDA.2017.8284101>.
- [10] BUDHI, G.S., CHIONG, R., PRANATA, I., and HU, Z., 2021. Using machine learning to predict the sentiment of online reviews: A new framework for comparative analysis. *Archives of Computational Methods in Engineering* 28, 2543-2566. DOI= <http://dx.doi.org/https://doi.org/10.1007/s11831-020-09464-8>.
- [11] BUDHI, G.S., CHIONG, R., and WANG, Z., 2021. Resampling imbalanced data to detect fake reviews using machine learning classifiers and textual-based features. *Multimedia Tools and Applications* 80, 13079-13097. DOI= <http://dx.doi.org/https://doi.org/10.1007/s11042-020-10299-5>.
- [12] BUDHI, G.S., CHIONG, R., WANG, Z., and DHAKAL, S., 2021. Using a hybrid content-based and behaviour-based featuring approach in a parallel environment to detect fake reviews. *Electronic Commerce Research and Applications* 47, 101048. DOI= <http://dx.doi.org/https://doi.org/10.1016/j.elerap.2021.101048>.
- [13] CAMPBELL, C. and YING, Y., 2011. *Learning with Support Vector Machines*. Morgan & Claypool.
- [14] CARDOSO, E.F., SILVA, R.M., and ALMEIDA, T.A., 2018. Towards automatic filtering of fake reviews. *Neurocomputing* 309, 106-116. DOI= <http://dx.doi.org/10.1016/j.neucom.2018.04.074>.
- [15] CHANG, C.C. and LIN, C.J., 2011. LIBSVM: A library for Support Vector Machines. *ACM Transactions on Intelligent Systems and Technology* 2, 3, 1-27. DOI= <http://dx.doi.org/https://doi.org/10.1145/1961189.1961199>.
- [16] CHIONG, R., BUDHI, G.S., and DHAKAL, S., 2021. Combining sentiment lexicons and content-based features for depression detection. *IEEE Intelligent Systems* 36, 6. DOI= <http://dx.doi.org/https://doi.org/10.1109/MIS.2021.3093660>.
- [17] CHIONG, R., BUDHI, G.S., DHAKAL, S., and CHIONG, F., 2021. A textual-based featuring approach for depression detection using machine learning classifiers and social media texts. *Computers in Biology and Medicine* 135, 104499. DOI= <http://dx.doi.org/https://doi.org/10.1016/j.combiomed.2021.104499>.
- [18] DENG, X., LI, Y., WENG, J., and ZHANG, J., 2019. Feature selection for text classification: A review. *Multimedia Tools and Applications* 78, 3, 3797-3816. DOI= <http://dx.doi.org/https://doi.org/10.1007/s11042-018-6083-5>.
- [19] DONG, M., YAO, L., WANG, X., BENATALLAH, B., HUANG, C., and NING, X., 2018. Opinion fraud detection via neural autoencoder decision forest. *Pattern Recognition Letters*. DOI= <http://dx.doi.org/10.1016/j.patrec.2018.07.013>.
- [20] ELMOGY, A.M., TARIQ, U., and MOHAMMED, A.I.A., 2021. Fake Reviews Detection using Supervised Machine Learning. *International Journal of Advanced Computer Science and Applications* 12, 1, 601-. DOI= <http://dx.doi.org/https://dx.doi.org/10.14569/IJACSA.2021.0120169>.

- [21] ETAIWI, W. and NAYMAT, G., 2017. The impact of applying different preprocessing steps on review spam detection. *Procedia Computer Science* 113, 273-279.
- [22] FELBERMAYR, A. and NANOPOULOS, A., 2016. The role of emotions for the perceived usefulness in online customer reviews. *Journal of Interactive Marketing* 36, 60-76. DOI= <http://dx.doi.org/10.1016/j.intmar.2016.05.004>.
- [23] FENG, V.W. and HIRST, G., 2013. Detecting deceptive opinions with profile compatibility. In *Proceedings of International Joint Conference on Natural Language Processing*, Nagoya, Japan, 338-346.
- [24] FREUND, Y. and SCHAPIRE, R.E., 1997. A Decision-Theoretic Generalisation of On-Line Learning and an Application to Boosting. *Journal of Computer and System Sciences* 55, 1 (1997/08/01/), 119-139. DOI= <http://dx.doi.org/https://doi.org/10.1006/jcss.1997.1504>.
- [25] FRIEDMAN, J.H., 2001. Greedy function approximation: A gradient boosting machine. *The Annals of Statistics* 29, 5, 1189-1232.
- [26] HAZIM, M., ANUAR, N.B., AB RAZAK, M.F., and ABDULLAH, N.A., 2018. Detecting opinion spams through supervised boosting approach. *PLoS ONE* 13, 6, e0198884. DOI= <http://dx.doi.org/10.1371/journal.pone.0198884>.
- [27] HERNÁNDEZ FUSILIER, D., MONTES-Y-GÓMEZ, M., ROSSO, P., and GUZMÁN CABRERA, R., 2015. Detecting positive and negative deceptive opinions using PU-learning. *Information Processing & Management* 51, 4, 433-443. DOI= <http://dx.doi.org/10.1016/j.ipm.2014.11.001>.
- [28] HEYDARI, A., TAVAKOLI, M., and SALIM, N., 2016. Detection of fake opinions using time series. *Expert Systems with Applications* 58, 83-92. DOI= <http://dx.doi.org/10.1016/j.eswa.2016.03.020>.
- [29] HEYDARI, A., TAVAKOLI, M.A., SALIM, N., and HEYDARI, Z., 2015. Detection of review spam: A survey. *Expert Systems with Applications* 42, 7, 3634-3642. DOI= <http://dx.doi.org/10.1016/j.eswa.2014.12.029>.
- [30] HU, Z., CHIONG, R., PRANATA, I., BAO, Y., and LIN, Y., 2019. Malicious web domain identification using online credibility and performance data by considering the class imbalance issue. *Industrial Management & Data Systems* 119, 3, 676-696. DOI= <http://dx.doi.org/https://doi.org/10.1108/IMDS-02-2018-0072>.
- [31] CHIONG, R., WANG, Z., FAN, Z., and DHAKAL, S., 2022. A fuzzy-based ensemble model for improving malicious web domain identification. *Expert Systems with Applications* 204, 117243. DOI= <https://doi.org/10.1016/j.eswa.2022.117243>.
- [32] JACOBS, R.A., JORDAN, M.I., NOWLAN, S.J., and HINTON, G.E., 1991. Adaptive Mixtures of Local Experts. *Neural Computation* 3, 1, 79-87. DOI= <http://dx.doi.org/https://doi.org/10.1162/neco.1991.3.1.79>.
- [33] JAVED, M.S., MAJEED, H., MUJTABA, H., and BEG, M.O., 2021. Fake reviews classification using deep learning ensemble of shallow convolutions. *Journal of Computational Social Science* 4, 2, 883-902. DOI= <http://dx.doi.org/10.1007/s42001-021-00114-y>.
- [34] CHIONG, R., FAN, Z., HU, Z., and DHAKAL, S., 2022. A novel ensemble learning approach for stock market prediction based on sentiment analysis and the sliding window method. *IEEE Transactions on Computational Social Systems*. DOI=<http://dx.doi.org/10.1109/TCSS.2022.3182375>.
- [35] JOHNSON, R.W., 2001. An introduction to the Bootstrap. *Teaching Statistics* 23, 2, 49-54. DOI= <http://dx.doi.org/https://doi.org/10.1111/1467-9639.00050>.
- [36] KERAS, 2019. Keras: The Python Deep Learning library.
- [37] KONTSCHIEDER, P., FITTERAU, M., CRIMINISI, A., and BUL'O, S.R., 2015. Deep Neural Decision Forests. In *2015 IEEE International Conference on Computer Vision (ICCV)* IEEE, Santiago, Chile, 1467-1475.
- [38] KUMAR, A., GOPAL, R.D., SHANKAR, R., and TAN, K.H., 2022. Fraudulent review detection model focusing on emotional expressions and explicit aspects: investigating the potential of feature engineering. *Decision Support Systems* 155. DOI= <http://dx.doi.org/10.1016/j.dss.2021.113728>.
- [39] KUMAR, N., VENUGOPAL, D., QIU, L., and KUMAR, S., 2018. Detecting review manipulation on online platforms with hierarchical supervised learning. *Journal of Management Information Systems* 35, 1, 350-380. DOI= <http://dx.doi.org/10.1080/07421222.2018.1440758>.
- [40] LECUN, Y., BOTTOU, L., BENGIO, Y., and HAFNER, P., 1998. Gradient-based learning applied to document recognition. *Proceedings of the IEEE* 86, 11, 2278-2324. DOI= <http://dx.doi.org/https://doi.org/10.1109/5.726791>.
- [41] LI, J., OTT, M., CARDIE, C., and HOVY, E., 2014. Towards a general rule for identifying deceptive opinion spam. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics*, Baltimore, Maryland, USA, 1566-1576.
- [42] LI, L., QIN, B., REN, W., and LIU, T., 2017. Document representation and feature combination for deceptive spam review detection. *Neurocomputing* 254, 33-41. DOI= <http://dx.doi.org/10.1016/j.neucom.2016.10.080>.
- [43] MALBON, J., 2013. Taking fake online consumer reviews seriously. *Journal of Consumer Policy* 36, 2, 139-157. DOI= <http://dx.doi.org/10.1007/s10603-012-9216-7>.
- [44] MANNING, C.D., RAGHAVAN, P., and SCHUETZE, H., 2008. Naïve Bayes text classification. In *Introduction to Information Retrieval* Cambridge University Press, 234-265.
- [45] MARTENS, D. and MAALEJ, W., 2019. Towards understanding and detecting fake reviews in app stores. *Empirical Software Engineering* 24, 6, 3316-3355. DOI= <http://dx.doi.org/10.1007/s10664-019-09706-9>.
- [46] MENARD, S., 2010. *Logistic Regression: From Introductory to Advanced Concepts and Applications*. SAGE, Los Angeles.
- [47] MOHAMMADI, Y., SALARPOUR, A., and CHOUHY LEBORGNE, R., 2021. Comprehensive strategy for classification of voltage sags source location using optimal feature selection applied to support vector machine and ensemble techniques. *International Journal of Electrical Power & Energy Systems* 124. DOI= <http://dx.doi.org/10.1016/j.ijepes.2020.106363>.
- [48] NORVIG, P., 2016. How to write a spelling corrector.
- [49] OTT, M., CARDIE, C., and HANCOCK, J.T., 2013. Negative deceptive opinion spam. In *Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, Atlanta, Georgia, US, 497-501.
- [50] OTT, M., CHOI, Y., CARDIE, C., and HANCOCK, J.T., 2011. Finding deceptive opinion spam by any stretch of the imagination. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics*, Portland, Oregon, 309-319.
- [51] PEDREGOSA, F., VAROQUAUX, G., GRAMFORT, A., MICHEL, V., THIRION, B., GRISEL, O., BLONDEL, M., PRETTENHOFER, P., WEISS, R., DUBOURG, V., VANDERPLAS, J., PASSOS, A., COUNAPEAU, D., BRUCHER, M., PERROT, M., and DUCHESNAY, É., 2011. Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research* 12, 85 (May 08), 2825-2830.
- [52] POLIKAR, R., 2006. Ensemble based systems in decision making. *IEEE Circuits and Systems Magazine* 6, 3, 21-45. DOI= <http://dx.doi.org/10.1109/MCAS.2006.1688199>.
- [53] QUINLAN, J.R., 1986. Induction of decision trees. *Machine Learning* 1, 1 (March 01), 81-106. DOI= <http://dx.doi.org/https://doi.org/10.1007/bf00116251>.
- [54] REN, Q., LI, M., and HAN, S., 2019. Tectonic discrimination of olivine in basalt using data mining techniques based on major elements: a comparative study from multiple perspectives. *Big Earth Data* 3, 1, 8-25. DOI= <http://dx.doi.org/https://doi.org/10.1080/20964471.2019.1572452>.

- [55] REN, Y. and JI, D., 2017. Neural networks for deceptive opinion spam detection: An empirical study. *Information Sciences* 385-386, 213-224. DOI= <http://dx.doi.org/10.1016/j.ins.2017.01.015>.
- [56] REN, Y., ZHANG, L., and SUGANTHAN, P.N., 2016. Ensemble Classification and Regression-Recent Developments, Applications and Future Directions [Review Article]. *IEEE Computational Intelligence Magazine* 11, 1, 41-53. DOI= <http://dx.doi.org/10.1109/mci.2015.2471235>.
- [57] ROUT, J.K., SINGH, S., JENA, S.K., and BAKSHI, S., 2016. Deceptive review detection using labeled and unlabeled data. *Multimedia Tools and Applications* 76, 3, 3187-3211. DOI= <http://dx.doi.org/10.1007/s11042-016-3819-y>.
- [58] RUMELHART, D.E., HINTON, G.E., and WILLIAMS, R.J., 1986. Learning internal representations by error propagation. In *Parallel Distributed Processing: Explorations in the Microstructure of Cognition* MIT Press, 318-362.
- [59] SAVAGE, D., ZHANG, X., YU, X., CHOU, P., and WANG, Q., 2015. Detection of opinion spam based on anomalous rating deviation. *Expert Systems with Applications* 42, 22, 8650-8657. DOI= <http://dx.doi.org/10.1016/j.eswa.2015.07.019>.
- [60] SCHAPIRE, R.E., 1990. The strength of weak learnability. *Machine Learning* 5, 2 (1990/06/01), 197-227. DOI= <http://dx.doi.org/10.1007/BF00116037>.
- [61] SHAN, G., ZHOU, L., and ZHANG, D., 2021. From conflicts and confusion to doubts: Examining review inconsistency for fake review detection. *Decision Support Systems* 144. DOI= <http://dx.doi.org/10.1016/j.dss.2021.113513>.
- [62] SHIN, J., YOON, S., KIM, Y., KIM, T., GO, B., and CHA, Y., 2021. Effects of class imbalance on resampling and ensemble learning for improved prediction of cyanobacteria blooms. *Ecological Informatics* 61. DOI= <http://dx.doi.org/10.1016/j.ecoinf.2020.101202>.
- [63] SUN, C., DU, Q., and TIAN, G., 2016. Exploiting product related review features for fake review detection. *Mathematical Problems in Engineering* 2016, 1-7. DOI= <http://dx.doi.org/10.1155/2016/4935792>.
- [64] TANG, X., QIAN, T., and YOU, Z., 2020. Generating behavior features for cold-start spam review detection with adversarial learning. *Information Sciences* 526, 274-288. DOI= <http://dx.doi.org/10.1016/j.ins.2020.03.063>.
- [65] WAHYUNI, E.D. and DJUNAIDY, A., 2016. Fake review detection from a product review using modified method of iterative computation framework. In *Proceedings of MATEC Web of Conferences* 58, 03003. DOI= <http://dx.doi.org/10.1051/matec>.
- [66] WANG, X., HE, K.L.S., and ZHAO, J., 2016. Learning to Represent Review with Tensor Decomposition for Spam Detection. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, Texas, US, 866-875.
- [67] WOLPERT, D.H., 1992. Stacked generalisation. *Neural Networks* 5, 2 (1992/01/01/), 241-259. DOI= [http://dx.doi.org/https://doi.org/10.1016/S0893-6080\(05\)80023-1](http://dx.doi.org/https://doi.org/10.1016/S0893-6080(05)80023-1).
- [68] WU, Y., NGAI, E.W.T., WU, P., and WU, C., 2020. Fake online reviews: Literature review, synthesis, and directions for future research. *Decision Support Systems* 132. DOI= <http://dx.doi.org/10.1016/j.dss.2020.113280>.
- [69] XIE, F., FAN, H., LI, Y., JIANG, Z., MENG, R., and BOVIK, A., 2017. Melanoma Classification on Dermoscopy Images Using a Neural Network Ensemble Model. *IEEE Trans Med Imaging* 36, 3 (Mar), 849-858. DOI= <http://dx.doi.org/10.1109/TMI.2016.2633551>.
- [70] XIONG, T., BAO, Y., HU, Z., and CHIONG, R., 2015. Forecasting interval time series using a fully complex-valued RBF neural network with DPSO and PSO algorithms. *Information Sciences* 305, 77-92. DOI= <http://dx.doi.org/https://doi.org/10.1016/j.ins.2015.01.029>.
- [71] ZHANG, D., ZHOU, L., KEHOE, J.L., and KILIC, I.Y., 2016. What online reviewer behaviors really matter? Effects of verbal and nonverbal behaviors on detection of fake online reviews. *Journal of Management Information Systems* 33, 2, 456-481. DOI= <http://dx.doi.org/https://doi.org/10.1080/07421222.2016.1205907>.
- [72] ZHANG, W., DU, Y., YOSHIDA, T., and WANG, Q., 2018. DRI-RCNN: An approach to deceptive review identification using recurrent convolutional neural network. *Information Processing & Management* 54, 4, 576-592. DOI= <http://dx.doi.org/10.1016/j.ipm.2018.03.007>.