

An Analysis of Machine Learning and Deep Learning Techniques for Detecting Sarcasm in Online Text

Student Name: Molly Hayward

Supervisor Name: Dr Noura Al-Moubayed

Submitted as part of the degree of BSc Computer Science to the
Board of Examiners in the Department of Computer Sciences, Durham University

Abstract —

Context / Background – Sarcasm is a unique form of language that cannot be interpreted at surface-level. Detecting sarcasm proves a significant challenge for traditional sentiment analysers and humans alike. This highlights the scope for innovative machine learning and deep learning solutions to this ill-structured problem.

Aims – Despite the challenges in this domain, this ultimate aim of this project is to produce a tool that can detect sarcasm in online text with a *high degree* of accuracy.

Method – State-of-the-art contextualised word embeddings (ElMo vectors) and deep learning classification models are used to realise this aim.

Results – Through extensive experimentation, we show that

Conclusions – In summary, we conclude that

Keywords — Machine learning, Deep learning, Sarcasm Detection, Sentiment, Classification

INTRODUCTION

Sarcasm is a complex linguistic phenomenon, which when present in text indicates that the literal interpretation differs from the implied meaning. Sarcasm is prevalent in online user-generated content, and poses a significant challenge within the field of natural language processing, specifically within sentiment analysis, as it inverts the sentiment-polarity of a positive or negative utterance. Sentiment analysis is the task of deducing the sentiment polarity of text - typically whether the author is in favour of, or against, a certain subject. It is increasingly common for organisations to use sentiment analysis in order to gauge public opinion on their products and services; however, classic sentiment analysers cannot deduce the implicit meaning of sarcastic text and will wrongly classify the author's opinion. Hence, any tool that strives to accurately determine the meaning of user-generated text must be capable of detecting sarcasm. The phenomenon of sentiment incongruity often lies at the heart of sarcasm, therefore advancements in automatic sarcasm detection research have the potential to vastly improve the sentiment analysis task.

A Problem Background

As sarcasm is multi-faceted, this makes its detection a unique and interesting challenge. It can be both explicit and implicit; oftentimes, contextual cues are more powerful indicators of sarcasm than the words themselves. However we often lose this context, and hence the sarcastic undertones, in transcription. Consider the scenario where a person is congratulated on their hard work,

despite the obviousness of them having not worked hard at all. Or perhaps a customer thanking a waiter for the delicious food, even though they sent it back to the kitchen. In isolation, the speech is not enough to convey the sarcastic intent. Furthermore, even humans struggle to consistently recognise sarcastic intent; due in part to the lack of an all-encompassing, universal definition of sarcasm. In a 2011 study, González-Ibáñez et al. [5] found low agreement rates between human annotators when classifying statements as either sarcastic or non-sarcastic, and in their second study, three annotators unanimously agreed on a label less than 72% of the time. Truly, the existence of sarcasm can only be conclusively affirmed by the author. Additionally, developmental differences such as autism, as well as cultural and societal nuances cause variations in the way different people define and perceive sarcasm.

These factors make sarcasm detection an extremely complex task for both humans and computers. Despite this, *most* humans can recognise sarcastic intent *most* of the time. We strive to replicate this performance, or even perhaps improve upon it, in order to move towards a more concrete definition of what *makes* a statement sarcastic. This highlights the scope for automating this process, and the need for an innovative solution.

Early rule-based approaches introduced the idea that the sentiment contained within text can be indicative of sarcasm, although exploring these sentiment features using more complex machine learning and deep learning architectures is a novel approach to the problem.

B Research Questions and Objectives

The **research questions** guiding this project are as follows –

- *Do deep learning techniques perform better than machine learning approaches?*
- *Can we improve the solution by incorporating sentiment features?*
- *Can a custom model identify which linguistic cues correlate more to sarcastic labels?*

The following objectives are designed to address these specific research questions, and are divided into three categories depending on their priority level and difficulty.

The **minimum** objectives of this project are to evaluate, compare and clean high-quality datasets, as well as to employ machine learning models and evaluate the effectiveness of these architectures on our chosen datasets. We found that a logistic regression model trained on ElMo vectors performed best for all three datasets, achieving F_1 scores of 84.6%, 81.3% and _%.

The **intermediate** objectives of this project are to experiment with and evaluate deep learning models on the chosen datasets, and to determine which is the best performing solution. Additionally, we evaluate the model against unseen data. This is achieved by ...

The **advanced** objectives of this project are to implement an attention-based deep neural network in order to identify words that correlate more to sarcastic labels, in order to produce a visualisation of attention words. We implement a custom attention layer using keras, and apply our trained model to a number of sentences to illustrate its effectiveness.

This is achieved by...

RELATED WORK

In this section, we compare classification model architectures used specifically in sarcasm detection research, as well as techniques that have seen success in other natural language processing (NLP) classification tasks. In most approaches, sarcasm detection is treated as a binary classification problem, in which text is grouped into two categories - sarcastic and non-sarcastic. It is otherwise treated as a multi-class problem where the extent to which a statement is sarcastic is ranked on a discrete scale, e.g. 1 to 5. In similar studies, the terms irony and sarcasm are often used interchangeably [23]. For clarity, an overview of their definitions are given in figure 1. Both definitions capture the humorous intent in both sarcasm and irony, however sarcasm extends this definition to include the potential for criticism - often present in opinion-based content such as product reviews.

Figure 1: Comparison of definitions of Irony and Sarcasm

Term	Definition
Irony	The use of words that are the opposite of what you mean, as a way of being funny. ¹
Sarcasm	The use of remarks that clearly mean the opposite of what they say, made in order to hurt someone’s feelings or to criticize something in a humorous way. ²

A Traditional Classifiers

A number of simple linguistic approaches have been used in previous research. One such class of naïve approaches is *rule-based*, where text is classified based on a set of linguistic rules. Maynard and Greenwood (2014) [12] proposed that the sentiment contained in hashtags can indicate the presence of sarcasm in tweets, such as #notreally in the example "I love doing the washing-up #notreally". They used a rule-based approach, whereby if the sentiment of a tweet contradicts that of the hashtag, it is labelled sarcastic. They reported an F_1 score of 91.03% on a random sample consisting of 400 tweets.

Riloff et al. (2013) [20] sarcasm is described as a contrast between positive sentiment and negative situation. This description is leveraged by Bharti et al (2015) [2], which presents two rule-based approaches to sarcasm detection. The first approach identifies sentiment bearing situation phrases, classifying the phrase as sarcastic if the negative situation is juxtaposed with a positive sentiment phrase. Rule-based techniques are fairly primitive when compared to their modern, deep learning counterparts, as each ruleset must be generated manually and for each dataset.

B Machine-Learning Classifiers

A machine-learning algorithm can learn patterns that are associated with sarcastic and non-sarcastic data. This allows it to classify unseen text into these categories without the need for a static linguistic rule set. If the training data is high-quality and plentiful, or an effective unsupervised-learning algorithm is used, then the model can begin to make accurate predictions. There are a number of commonly used architectures of machine-learning classifiers, including Naïve Bayes, Support Vector machines (SVMs) and logistic regression based classifiers. Although, not all approaches fit neatly into these categories.

¹www.dictionary.cambridge.org/dictionary/english/irony

²www.dictionary.cambridge.org/dictionary/english/sarcasm

Tsur et al. (2010) [23] proposed a novel semi-supervised algorithm, *SASI*, for sarcasm identification; trained on a small balanced seed of 160 reviews (80 sarcastic and 80 non-sarcastic) from a corpus of 66000 human-annotated Amazon product reviews. In order to mimic the ambiguous and spectral nature of sarcasm, a discrete score between 1 (non-sarcastic) and 5 (definitely sarcastic) is assigned to each sentence; scores of 3 and higher indicate sarcasm. Syntactic and pattern-based features are extracted and fed to a classifier which utilises a strategy similar to k-nearest neighbor (kNN), whereby a sarcasm score is assigned to an unseen vector based upon the weighted average of the k-nearest neighbors in the training set. They reported an F_1 score of 82.7% on the binary classification task

Naïve Bayes classifiers are supervised algorithms that use bayes theorem to make predictions. Reyes et al. (2013) [19] used a naïve bayes and decision trees algorithm in order to detect irony which is closely related to sarcasm. They experimented with balanced and unbalanced distributions (25% ironic tweets and 75% other), achieving an F-score of 0.72 on the balanced distribution, dropping to 0.53 for the imbalanced distribution. In a similar vein, Barbieri et al. (2014) [1] used random forest and decision tree classifiers, also for the detection of irony. They used six types of features in order to represent tweets, and recorded results over three categories of training data - education, humor and politics.

Support Vector Machines (SVM), can be effective when smaller amounts of training data are available, however they require more computational resources than naïve bayes. SVMs use kernel functions to draw hyperplanes dividing a space into subspaces. González-Ibáñez et al. (2011) [5] used two classifiers, support vector machine with sequential minimal optimisation, as well as logistic regression. Their best result was an accuracy of 0.65, achieved with the combination of support vector machine with sequential minimal optimisation and unigrams. Ptáček et al. (2014) [18] also used support vector machine and Maximum Entropy classifiers. They also performed classification on balanced and imbalanced distributions.

C Deep-Learning Classifiers

Deep neural networks (DNNs) are increasingly being used in text classification tasks [17, 26]. Though there are many forms of DNN, two commonly used architectures are Recurrent Neural Networks (RNNs) and Convolutional Neural Networks (CNNs). Oftentimes, deep neural networks require a lot more training data than machine learning algorithms, but have been shown to produce state-of-the-art results in several domains.

One particularly interesting technique is the *Hierarchical Attention Network (HAN)* defined in Yang et al. (2016) [24] - an attention-based LSTM. The intuition is that certain words and sentences in a document contribute more to its overall meaning, and this is highly context dependent. They included one attention mechanism at word-level and another at sentence-level, and this allowed them to determine which words and phrases correlate more to certain labels. Using a similar approach in the domain of sarcasm may allow me to highlight the attention-words that are strongly influence the level of sarcasm. A similar technique was applied in Ghosh et al. (2018) [4] where they concluded that emoticons such as ':' and interjections e.g. "ah", "hmm" correspond with highly weighted sentences. Gated Recurrent Units (GRUs) are another type of RNN used to overcome the vanishing gradient problem. Zhang et al. (2016) [25] used a bidirectional gated RNN to capture local information in tweets, as well as a pooling neural network to extract context from historic tweets, achieving 78.55% accuracy.

Convolutional Neural Networks allow us to extract higher-level features, and they are based on the mathematical operation of convolution. Zhang et al. (2015) [26] explored the used of

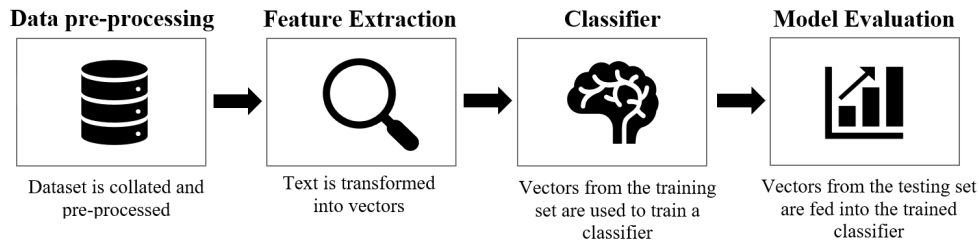
character-level convolutional neural networks for text classification. This network consists of 6 convolutional layers and 3 fully-connected layers. They found that it showed better performance on raw texts such as an Amazon product review corpus. CNNs have been previously used in sarcasm detection. For example, *Poria et al. (2017)* [17] describes a convolutional neural network (CNN) for sarcasm detection.

SOLUTION

This project addresses a binary classification problem, aiming to determine if an unseen snippet of text is sarcastic or non-sarcastic. Our sarcasm detection tool employs trained classification models which form predictions based on the knowledge gained from features encountered during training. Developing highly accurate classification models remains the greatest challenge in this project, therefore it is the central focus of our solution. Model performance is strongly dependent on the data pre-processing and feature extraction techniques employed, as well as the choice of model architecture and training procedures.

In our solution, high quality datasets are first collected and pre-processed. During feature extraction, the textual training data is transformed into numerical vectors and these features are used to train a classification model. The trained classifier is then carefully evaluated in order to account for potential overfitting. The following figure outlines the model pipeline, and each of these stages are further addressed in the remaining sections.

Figure 2: Overview of the model pipeline



A Implementation Tools

This project is majoritively written in **Python**, as it has an abundance of open-source libraries well-suited to machine learning and deep-learning tasks. Additionally, Python is very readable, making it easy for others to interpret, test and reproduce code. We use **pandas** for manipulating the datasets, as it provides functionality for applying functions globally across the dataset. We use **pickle** to store the vectorised datasets in memory. Vectorisation is a costly process, therefore writing the files to memory reduces long-term time demands during the experimentation and testing phases. We also incorporate **SpaCy** when cleaning and tokenising the data.

In terms of *classification*, we implement most of the machine learning architectures using **scikit-learn** - an open-source machine learning library. For constructing deep learning models, we use **keras**, as its modular structure allows for extensive experimentation and simplifies tweaking of model architectures.

B Datasets

In similar studies, large datasets have been collected and some have been made publically available. We experiment with data from multiple domains and online platforms, including Twitter, Amazon, and Online News outlets, giving broad coverage of sarcasm usage in formal and informal media. The suitability of a dataset is determined based on size, cleanliness and composition.

Small datasets do not provide a sufficiently representative view of the true distribution of data, therefore a model trained on this data will fail to generalise to unseen data. On the contrary, it may be infeasible to train a model in a reasonable length of time using large datasets as this will require vast amounts of computational resources. For this reason, we use datasets with fewer than 50000 instances.

Secondly, human-annotated data with minimal spelling and grammatical errors is preferable to data that is automatically selected using rule-based heuristics (e.g. tweets annotated with #sarcasm) as it is typically of higher quality, however it is not always possible to gather this data in large quantities due to the time overhead involved. Additionally, we prefer to use balanced datasets; despite sarcasm being a minority class due to its infrequent usage in real-life conversation. The reason for this is illustrated in the following example: if we train a classifier on an unbalanced binary dataset whereby the class distribution is 95% to 5%, we can achieve an accuracy of 95% simply by annotating each datapoint with the majority class label. Most machine learning algorithms are designed to optimise overall accuracy, therefore it is important to provide it with approximately equal amounts of sarcastic and non-sarcastic data.

Three datasets are selected, and a statistical break down of their properties is provided in *Figure 3*, including the number of instances in the dataset, the average number of tokens per instance and the percentage of tokens that are included in the dictionary. GloVe is a method of feature extraction where word embeddings for tokens are provided as a dictionary. The number of tokens present in the GloVe dictionary is a measure of quality, as this method of feature extraction cannot be used on data containing many ill-formed and unrepresented tokens. All tokens are converted to lowercase beforehand to make for a fair comparison between datasets, as only lowercase tokens are present in the GloVe dictionary. A detailed explanation of GloVe is provided in the feature extraction section. A benchmark score is provided where the dataset was used in the original paper to train one or more classifiers. The best score achieved in each paper is used as the benchmark which we aim to exceed.

Figure 3: Statistical breakdown of the datasets used in this project
N.B. *pos* indicates positive instances and *neg* indicates negative instances.

Dataset	Number of data points	Avg. tokens	% in GloVe dictionary	Benchmark Score
News Headlines Dataset For Sarcasm Detection	28619 pos: 47.6% neg: 52.4%	11.2	97.8%	0.897
Sarcasm Amazon Review Corpus	1254 pos: 34.9% neg: 65.1%	276.4	98.2%	N/A
Ptáček et al. (2014)	17280 pos: 50.4% neg: 49.6%	19.8	98.2%	0.947

1. News Headlines Dataset For Sarcasm Detection

This dataset was collated for use in *Misra et al. (2019)* [14]. Headlines are collected from two news sources - *The Onion* ¹ and *The Huffington Post* ². The Onion is renowned for posting satirical news, and these sarcastic headlines account for the sarcastic instances in the dataset. The non-sarcastic instances were collected from the Huffington Post. There is little noise in this dataset, as news headlines are written formally, and are far less likely to contain spelling and grammar errors than data from more informal platforms. These headlines are self contained, unlike tweets which may be in response to another tweet, thereby excluding necessary context. Misra et al. (2019) [14] used a hybrid neural network architecture and achieved an *accuracy* of 0.897. NB. Accuracy was used as the dataset is approximately balanced.

2. Sarcasm Amazon Review Corpus

This corpus was collected using the procedure outlined in *Filatova (2012)* [3]. Amazon reviews are collected via crowdsourcing, and are human-labelled by 5 annotators. The number of stars given to a review, as well as the title and content are provided. These reviews are typically longer than the average tweet, which is capped at 280 characters. The author's punctuation and spelling is preserved, therefore this data is very messy.

3. Ptáček et al. (2014)

This dataset is one of three datasets collated in *Ptáček et al. (2014)* [18]. We use the balanced dataset of 100 000 English tweet ids, whereby the original tweet must be scraped using the Twitter API. Of this original corpus, a small portion of tweets (_ tweets) are no longer available on Twitter, therefore they are omitted from the dataset. This dataset is messy, as sarcastic instances are those that include *#sarcasm*, and non-sarcastic instances are selected arbitrarily from the set of all other tweets. Hence, some of the sarcastic instances may not be sarcastic, and some non-sarcastic instances may be sarcastic. We take a subset of the remaining tweets (7500 from each class) to form a smaller balanced corpus. Ptáček et al. (2014) [18] achieved an F_1 score of 0.947.

¹www.theonion.com

²www.huffpost.com

C Data pre-processing

Data pre-processing is an important stage in natural language processing, with the potential to hinder or boost model performance. Datasets in this domain are inherently messy, often because they are collated from user generated content. Hence, data pre-processing is vital to reduce noise and sparsity in the feature space caused by inconsistent letter capitalisation, varying use of punctuation and erroneous spelling. The suitability of each technique is highly dependent on the dataset and on the feature-extraction technique, as certain types of pre-processing may result in the removal of useful features.

Use of #not may, the quality of the dataset may be dubious.

Reducing sparsity in the dataset is vital in order to use word embeddings such as Word2Vec and GloVe. They cannot generalise to out-of-vocabulary tokens i.e. tokens that were not in the original training set, and in the GloVe dictionary, all tokens are given in lowercase. As a result, it is essential that the data is converted to lowercase, however capitalisation can indicate emphasis, which may be indicative of sarcasm. N.B. It is necessary to re-introduce this meaning in some form. We add the token ! before a capitalised word, indicating emphasis, or "shouting". On twitter, users can include hyperlinks to other websites in their tweets, and this can be a source of noise in the dataset. Ptáček *et al.* (2014) [18] collated a dataset of sarcastic and non-sarcastic tweets, whereby the presence of sarcasm is indicated by the marker #sarcasm, removing URLs and references to users, as well as hashtags. However, Liebrecht *et al.* (2013) [10] showed that data contained in hashtags can be used to indicate sarcasm, such as #not. Transforming text into lowercase is a common pre-processing strategy, useful for reducing sparsity caused by variations in letter capitalisation e.g. 'Cat', 'cAt', and 'CAT' are all mapped to the same word - 'cat'. Similarly, removing punctuation and extending contractions (e.g. don't → do not) is commonplace. However, punctuation and capitalisation can be used for emphasis, therefore removing these attributes may make the presence of sarcasm less obvious.

In the *News Headlines Dataset For Sarcasm Detection* dataset, the raw data was pre-cleaned, as it had all been transformed into lowercase. As news headlines are written in a formal manner, spelling and grammar errors are rare, thereby increasing the likelihood that pre-trained embeddings are present in the GloVe dictionary.

D Feature Extraction

Extracting meaningful data from large corpora is a complex task, however in order to do this successfully, we first construct word vectors that capture semantic and syntactic relationships. The textual data is transformed into numerical vectors, so as to encapsulate the seminal features of language.

Traditional Approaches – We utilise a number of traditional vectorisation methods - perhaps the simplest of which is *Bag of Words (BOW)*, whereby a snippet of text is encoded as a single vector containing a count of the number of occurrences of each word. *Term Frequency-Inverse document frequency (TF-IDF)* [21] extends the BOW model by providing each word with a score that reflects its level of importance based upon its frequency within the document. These approaches have two main disadvantages, firstly, each vector is the size of the vocabulary, creating sparse inefficient representations. Secondly, vectors produced in this way do not preserve the context

within which words are found, and sarcasm is a highly contextual phenomenon. Nonetheless, they act as a good baseline against which to compare more advanced vectorisation approaches.

Word Embedding –

Embeddings produce lower-dimensional vectors than the sparse representations generated by earlier techniques. Introduced in Mikolov et al. (2013a) [13], *Word2Vec* describes a group of related models that can be used to produce high-quality vector representations of words - two such models are *Continuous Bag-of-Words* and *Continuous Skip-Gram*. It consists of a shallow (two-layer) neural network that takes a large corpus and produces a high-dimensional vector space, such that words that share a similar context are clustered together in the feature space. Word2Vec is able to preserve the semantic relationships between words, even constructing analogies by composing vectors e.g. king - man + woman \approx queen. Likewise, it captures syntactic regularities such as the singular to plural relationship e.g. cars - car \approx apples - apple. In Word2Vec, intrinsic statistical properties of the corpus, which were key to earlier techniques, are neglected, therefore global patterns may be overlooked. To mitigate this, the *GloVe* [15] approach generates word embeddings by constructing an explicit word co-occurrence matrix for the entire corpus, such that the dot product of two word embeddings is equal to log of the number of times these words co-occur (within a defined window). An example is given in figure .. This allows it to capture both global and local statistics of the corpus. GloVe is used in this project to produce word embeddings for each corpus, as some of the datasets have a high average number of tokens in each snippet, therefore local context is insufficient.

Despite their ability to preserve semantic relationships, Word2Vec and GloVe do not accommodate polysemy, which describes the co-existence of alternative meanings for the same word. In addition, they cannot generalise to words that were not specifically included in the training set.

Figure -: Word co-occurrence matrix with window size 1 for the snippet:

"the clouds are in the sky"

	the	clouds	are	in	sky
the	0	1	0	1	1
clouds	1	0	1	0	0
are	0	1	0	1	0
in	1	0	1	0	0
sky	1	0	0	0	0

Given a corpus containing W tokens, a $W \times W$ matrix is constructed. The value in cell (x,y) denotes the frequency with which token x has co-occurred with token y within a fixed-size window. From this co-occurrence matrix, the probability of token x co-occurring with token y is given by the formula:

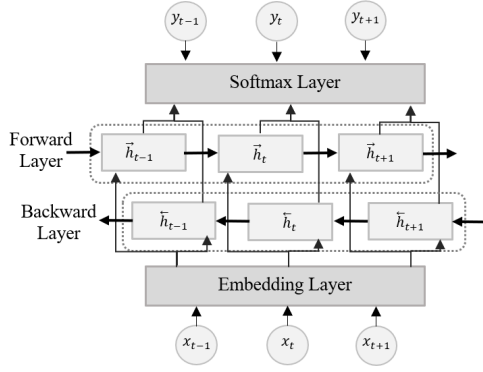
$$P_{co}(t_x|t_y) = \frac{C(t_y, t_x)}{C(t_y)}$$

where $C(t_i, t_j)$ is the number of times token i co-occurs with token j .

The word meanings are encapsulated by the ratios of co-occurrence probabilities, rather than the probabilities themselves.

Contextualised Word Embedding – Deep neural language models offer a more robust approach. Peters et al (2018) [16] introduced the *ELMo* model, and showed that the addition of ELMo to existing models can markedly improve the state-of-the-art in a number of NLP tasks.

Figure -: ELMo model pipeline



ELMo utilises a bi-directional LSTM (*long short-term memory*) model, concatenating the left-to-right and right-to-left LSTM in order to produce deep contextualised word representations. They are character based, therefore robust representations can be constructed for out-of-vocabulary tokens. However, rather than 'looking-up' pre-computed embeddings, ELMo generates them dynamically. Hence, it can disambiguate the sense of a word given the context in which it was found. These state-of-the-art contextualised embeddings are used in this project as the most advanced method of feature extraction.

Experimentation with miscellaneous features – As supervised machine learning requires hand-crafted features to inform predictions, we consider a number of hand-crafted features in addition to the vectors produced by the aforementioned vectorisation approaches. This includes lexical and sentiment features, as well as meta data where this is available as part of the original dataset. For example, the Sarcasm Amazon Review Corpus contains the number of stars in a review, and we can leverage this information as an additional feature. Features that can be shown to improve classification scores can be concatenated with our existing vectors (enrichment).

The following experiment is designed to answer the research question: *Can we improve the solution by incorporating sentiment features?* We use a state-of-the-art model, the SentimentAnnotator [22] by Stanford CoreNLP [11], in order to extract sentiment features from each of our datasets. The sentiment annotator assigns a discrete label between 0 (very negative) and 4 (very positive) to a snippet of text. For each instance in the dataset, we collect the sentiment values for each of the individual tokens, as well as the overall sentiment of the snippet. We construct 1×6 dimensional feature vectors, where values 0 - 4 show the distribution of sentiment classes for tokens in the snippet, and value 5 shows the overall sentiment label.

Example: "i love when my train is late!"

Tokens	i	love	when	my	train	is	late	!
Sentiment Labels	2	4	2	2	2	2	2	2

Sentiment Labels	0	1	2	3	4
Frequency of tokens	0	0	7	0	1
Proportion of tokens	0	0	0.875	0	0.125

Overall sentiment: 3.0

Feature vector: [0.0, 0.0, 0.875, 0.0, 0.125, 3.0]

E Classification

In this section, machine-learning and deep-learning classification models are considered. These algorithms are trained a set of rich features as outlined in the previous section.

E.1 Machine Learning Models

We experiment with five candidate machine learning classifiers, evaluating each of them when trained on four types of features. Each of these supervised models perform binary classification, where unseen samples are predicted to be sarcastic or non-sarcastic. A brief description of the models are provided below.

1. **K-Nearest Neighbours (KNN)** – The KNN classifier is a non-parametric model which uses lazy learning at the time of prediction. An unseen sample is assigned the most common class label amongst its k nearest neighbours, where k is a hyperparameter whose value is finetuned for each dataset.
2. **Support Vector Machine (SVM)** – A hyperplane is drawn such that the positive and negative features are partitioned into two groups. The features at the shortest euclidean distance to the hyperplane are known as support vectors, and their distance from the hyperplane is known as the margin. An SVM attempts to maximise this margin in order to create a decisive boundary between the sarcastic and non-sarcastic data.
3. **Naïve Bayes Classifier** – This is a probabilistic model derived from Bayes theorem, where all features are assumed to be independent of one another. Commonly used in the document classification problem, we calculate $P(\textit{Sarcastic}|\textit{features})$. This is the probability of an outcome A (i.e. statement is sarcastic), given its feature vector, x . There are a number of variations of Naïve Bayes - including Multinomial, Bernoulli and Gaussian.
4. **Logistic Regression** – A statistical model that aims to predict the likelihood of an event occurring given some previous data. It works with binary data, i.e. is the statement sarcastic or not. It is a more advanced approach than linear regression, which aims to model data using a linear equation.
5. **Random Forest Classifier** – An ensemble learning method, the random forest classifier fits a number of decision tree classifiers on subsamples of the training set. The decision trees are fairly independent of one another, and this low correlation allows them to produce ensemble predictions that outperform their individual predictions. Where one model goes wrong, the others are able to contradict this incorrect prediction.

E.2 Deep Learning Models

Deep learning classifiers do not need to use hand-crafted features used in supervised classifiers. We use Binary Cross-Entropy loss.

E.2.1 Recurrent Neural Networks – In a Recurrent Neural Network (RNN), prior inputs are used to inform future outputs. As in all neural networks, input vectors are modified by weighted nodes as they travel through the network. In a RNN, there is an additional hidden state which represents the context based on prior inputs, and is updated by inputs as they are fed through the network. This gives way to the interesting property that the same input could produce a different output depending on the previous inputs given to the RNN. In a vanilla RNN, the input and hidden state are simply propagated through a single tanh layer. However in a Long Short Term Memory model (LSTM), three gates are introduced, as well as a cell state, allowing an LSTM to better preserve long-term dependencies.

RNNs have been successfully applied to language-related tasks. Often, we may require more context than just the most recent surrounding text. The gap between the most relevant information needed to form a prediction can be very wide, if the relevant contextual information was given at the beginning of the paragraph. LSTMs are especially suited to processing sequential, or time-series, information such as text. This is due to the fact that RNNs have a 'memory', in that the input to the current step is the output from the previous step. However, they suffer from the vanishing gradient problem. This is where gradient values get smaller as it backtracks through lower layers, making it harder for the network to update weights and causing calculations to take longer. Recurrent Neural Networks are not very good at remembering long term dependencies, therefore in this instance we can use *Long short-term memory (LSTM)* models [7] instead which are more effective at preserving long term dependencies.

E.2.2 Convolutional Neural Networks – A Convolutional Neural Network (CNN) [9] is a specialised form of neural network. They were popularised for their applications in computer vision, hence they are designed to be used with two-dimensional image data [8], but have since been adapted for text-based classification tasks. CNNs consist of a few layers of convolutions with non-linear activation functions - they have fewer layers than typical neural networks, where each layer applies different filters and combines their result. The convolution operation is very fast, therefore many convolutional layers can be used in a single network.

Figure -: $m \times n$ feature vector

$$\begin{bmatrix} w_{1,1} & w_{1,2} & \dots & w_{1,N} \\ w_{2,1} & w_{2,2} & \dots & w_{2,N} \\ \vdots & \vdots & \vdots & \vdots \\ w_{M,1} & w_{M,2} & \dots & w_{M,N} \end{bmatrix}$$

Each token in a sentence is represented independently by a fixed-size n -dimensional vector; as such, the sentence is represented by an $m \times n$ vector of feature weights, where m is the number of tokens in the sentence. This feature vector is used as input to the convolutional neural network.

Convolutional layers compute the dot product of a kernel vector of k weights with k rows in the feature vector (representing a k -tokens in the input sentence), producing a new sequence of features. The convolution filter is repeatedly applied to overlapping sections of the matrix, and a max pooling operation is applied over each resulting feature map. Tokens in a sentence occur on a time axis, as the order that the tokens appear in gives rise to their meaning; hence, the filter slides downwards over the rows of the vector. Applying *maximum pooling over time* produces a single output value following each convolution, such that the feature corresponding to a particular vector is the maximum activation in the feature map. Producing a fixed-size output matrix means that variable size sentences and variable size filters can be used to obtain features

with the same output dimensions. This is especially useful for classification tasks, where the desired output is a fixed-sized vector indicating which class the input vector belongs to.

RESULTS

A Evaluation Method

This section discusses our approach to evaluating the successes and failures of the trained models. A simple approach is to use accuracy which refers to the proportion of data that is correctly labelled as either sarcastic or non-sarcastic. However, sarcasm is a minority class therefore on an unbalanced dataset we could achieve high accuracy by simply labelling every statement as non-sarcastic. In an attempt to mitigate against this, we instead form a conclusion based on the F_1 score i.e. the harmonic mean of precision and recall, where scores range from 0 (worst) to 1 (best). This metric has faced some criticism for giving equal weight to both precision and recall [6], therefore we consider both measures separately, as well as in combination.

$$\text{precision} = \frac{\text{true positives}}{\text{true positives} + \text{false positives}} \quad \text{recall} = \frac{\text{true positives}}{\text{true positives} + \text{false negatives}}$$

In the domain of sarcasm detection, precision refers to the proportion of the classified-sarcastic data that is *truly sarcastic* i.e. how many of the positives are true positives, and recall describes the proportion of the truly sarcastic data in the corpus that is *classified* as such i.e. how many true positives are labelled as positives.

B Feature Analysis

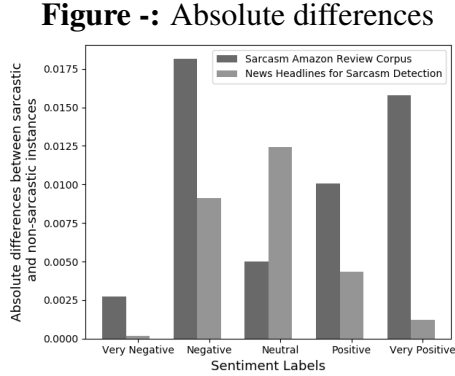
We analyse the effectiveness of the hand-crafted features, sentiment, punctuation and topic modelling, on each of the three datasets. As the feature vectors are small, we use logistic regression for this task. We train three logistic regressors (one for each dataset) separately, using each feature type in isolation, then use 5-fold cross-validation to compute the F_1 scores of the trained models. This enables us to assess whether these features alone are indicative of sarcasm in different media types. We deem scores above 0.5 to be indicative of useful features, particularly for scores significantly greater than this baseline.

Figure -: F_1 score of logistic regression models trained on hand-crafted features

Dataset	Feature Type		
	Sentiment	Punctuation	Topic
News Headlines for Sarcasm Detection	0.516	0.344	0.525
Sarcasm Amazon Review Corpus	0.685	0.394	0.520
Ptáček et al. (2014)	0.593	0.611	0.530

Our results show that sentiment features are indicative of sarcasm when extracted from the Sarcasm Amazon Review Corpus, however when extracted from the News Headlines for Sarcasm Detection, they were ineffective. To further examine this relationship between sarcasm and sentiment, we compute the mean proportion of tokens in each sentiment class for sarcastic and

non-sarcastic instances. We then compute the absolute differences of these values between the sarcastic and non-sarcastic classes - the intuition being that it is simpler for a classifier to separate samples where there is a clearer distinction between the expected features of each class. Hence, where the absolute difference in the mean proportion of tokens between sarcasm classes is higher, a classifier is better able to correctly annotate samples. A visualisation of these results is given in the figure below.



For the Sarcasm Amazon Review Corpus, within most sentiment classes there is a greater difference in the proportion of tokens for non-sarcastic and sarcastic instances. This indicates that sarcastic and non-sarcastic instances are easier to distinguish, given that their sentiment compositions are more distinct. The News Headlines for Sarcasm Detection results displays the opposite, and this explains why the logistic regression classifier trained on this data achieved an F_1 score of only 0.515, as the features for each class are very similar.

C Machine Learning Experimentation

We train a range of classifiers independently on each of the three datasets, and provide the results of this experimentation in the following subsections. The F_1 scores of the trained models are provided to 3 significant figures, and the score of the best-performing model is highlighted for clarity. Each classifier is trained on four types of vectors, Bag of Words, TF-IDF, GloVe and ElMo as outlined in the solution.

We also experimented with an unsupervised model - K-Means clustering where the majority label of each cluster was assigned, however, the results were very poor.

Experimental Settings – For the machine learning experiments, 5-fold cross validation was employed to ensure reproducibility and consistency of results. We allocate 80% of data to the training set, and the remaining 20% to the testing set for each of the five folds. For the K-nearest neighbor classifier, $K = 5$ was found to be the greatest balance between computational cost and noise mitigation.

Figure -: F_1 scores of Machine Learning classifiers on the *News Headlines for Sarcasm Detection* dataset

Model Architecture	Vector			
	BOW	TF-IDF	GloVe	ElMo
Support Vector Machine	0.714	0.718	0.732	0.844
Logistic Regression	0.721	0.712	0.732	0.846
Random Forest Classifier	0.530	0.643	0.754	0.797
Gaussian Naive Bayes	0.662	0.649	0.703	0.640
K-Means	0.400	0.433	0.501	0.471

Model Architecture	F ₁ Score
Logistic Regression + ElMo + Sentiment	0.846
Logistic Regression + ElMo + Punctuation	0.847
Logistic Regression + ElMo + Topic	0.847
Logistic Regression + ElMo + Topic + Punctuation	0.847

Figure -: F_1 scores of Machine Learning classifiers on the *Sarcasm Amazon Review Corpus*

	Vector			
Model Architecture	BOW	TF-IDF	GloVe	ElMo
Support Vector Machine	0.736	0.775	0.700	0.800
Logistic Regression	0.745	0.748	0.713	0.809
Random Forest Classifier	0.525	0.487	0.669	0.727
Gaussian Naive Bayes	0.625	0.621	0.416	0.666
K-Means	0.474	0.334	0.353	0.316

Model Architecture	F ₁ Score
Logistic Regression + ElMo + Sentiment	0.830
Logistic Regression + ElMo + Punctuation	0.809
Logistic Regression + ElMo + Topic	0.811
Logistic Regression + ElMo + Sentiment + Topic	0.811

Figure -: F_1 scores of Machine Learning classifiers on the *Ptáček et al. (2014) dataset*

	Vector			
Model Architecture	BOW	TF-IDF	GloVe	ElMo
Support Vector Machine	0.713	0.724	0.635	0.754
Logistic Regression	0.742	0.748	0.630	0.762
Random Forest Classifier	0.723	0.715	0.645	0.735
Gaussian Naive Bayes	0.681	0.675	0.594	0.684
K-Means	0.378	0.493	0.350	0.338

Model Architecture	F ₁ Score
Logistic Regression + TF-IDF + Sentiment	0.763
Logistic Regression + TF-IDF + Punctuation	0.762
Logistic Regression + TF-IDF + Topic	0.763
Logistic Regression + TF-IDF + Sentiment + Topic	0.764

D Deep Learning Experimentation

Each deep learning model is trained separately on each of the three datasets, using glove vectors and elmo vectors. Sparse feature representations such as tf-idf and bag of words are not used as

D.1 Experimental Settings

The upper-bound on the number of epochs is 300, although this is seldom reached due to the use of the Early Stopping callback function. The patience is set to 20, and this is the number of epochs for which no improvement can be seen before the training is prematurely terminated. However, setting the patience too low can hamper a model that plateaus for a while before continuing to improve.

1. **Early stopping** – Training is interrupted when the validation loss stops declining, using thins to determine the appropriate number of epochs.
2. **Model Checkpointing** – We track and save the weights of the best performing model during training; the model with the lowest validation loss.

Preliminary experimentation showed that the training process on the Sarcasm Amazon Review corpus was very unstable as this dataset is small. Results varied widely, and model scores were significantly lower. We augment this particular dataset using a synonym replacement strategy, whereby for each instance of the dataset, a new instance is added. This is a modified version of the

Figure -: F_1 scores of Deep Learning classifiers on the *News Headlines for Sarcasm Detection* dataset

Model Architecture	Vector	
	GloVe	ELMo
CNN	0.823	
LSTM	0.865	0.
Bi-directional LSTM	0.864	0.
Vanilla RNN	0.817	0.
Vanilla GRU	0.852	0.

Figure -: F_1 scores of Deep Learning classifiers on the *Sarcasm Amazon Review Corpus*

Model Architecture	Vector	
	GloVe	ELMo
CNN	0.	0.
LSTM	0.	0.
Bi-directional LSTM	0.	0.

Figure -: F_1 scores of Deep Learning classifiers on the *Ptáček et al. (2014) dataset*

	Vector	
Model Architecture	GloVe	ElMo
CNN	0.	0.
LSTM	0.	0.
Bi-directional LSTM	0.	0.

EVALUATION

In the following section, we analyse the effectiveness of the solution, discussing its strengths and weaknesses and evaluating to what extent it satisfies the research questions and deliverables.

A Solution Strengths

We found that sentiment features can be useful predictors of sarcasm in *some* forms of media. This may be due to the fact that sarcasm must be more overt in product reviews in order to mock the product, rather than the subtle, more contextualised instances in News Headlines.

B Solution Limitations

It was only feasible to evaluate a small range of online media types, however surveying a broader range of media types could give us a wider perspective of sarcasm in online media.

C Lessons learnt

We learned that

D Project organisation and approach

CONCLUSIONS

In terms of satisfying the objectives, this project was an overall success. With an interactive console, it may be useful to individuals who struggle to detect sarcasm in their daily-life, or for the industrial applications of enhanced opinion mining.

The main contributions of this project are as follows: sentiment features are combined with state-of-the-art vectorisation approaches.

For future work, it would be interesting to explore the use of sarcasm in a wider variety of media types, and to extend support to additional non-English languages. Exploring the use of transfer learning with Transformers (BERT) may also produce interesting results.

Table 1: SUMMARY OF PAGE LENGTHS FOR SECTIONS

Section		Number of Pages
I.	Introduction	2–3
II.	Related Work	2–3
III.	Solution	4–7
IV.	Results	2–3
V.	Evaluation	1–2
VI.	Conclusions	1

References

- [1] Francesco Barbieri and Horacio Saggion. Modelling irony in twitter. In *Proceedings of the Student Research Workshop at the 14th Conference of the European Chapter of the Association for Computational Linguistics*, pages 56–64, 2014.
- [2] Santosh Kumar Bharti, Korra Sathya Babu, and Sanjay Kumar Jena. Parsing-based sarcasm sentiment recognition in twitter data. In *Proceedings of the 2015 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining 2015*, pages 1373–1380. ACM, 2015.
- [3] Elena Filatova. Irony and sarcasm: Corpus generation and analysis using crowdsourcing. In *Lrec*, pages 392–398. Citeseer, 2012.
- [4] Debanjan Ghosh, Alexander R Fabbri, and Smaranda Muresan. Sarcasm analysis using conversation context. *Computational Linguistics*, 44(4):755–792, 2018.
- [5] Roberto González-Ibáñez, Smaranda Muresan, and Nina Wacholder. Identifying sarcasm in twitter: a closer look. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies: Short Papers-Volume 2*, pages 581–586. Association for Computational Linguistics, 2011.
- [6] David Hand and Peter Christen. A note on using the f-measure for evaluating record linkage algorithms. *Statistics and Computing*, 28(3):539–547, 2018.
- [7] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- [8] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.
- [9] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- [10] CC Liebrecht, FA Kunneman, and APJ van Den Bosch. The perfect solution for detecting sarcasm in tweets# not. 2013.
- [11] Christopher D Manning, Mihai Surdeanu, John Bauer, Jenny Rose Finkel, Steven Bethard, and David McClosky. The stanford corenlp natural language processing toolkit. In *Proceedings of 52nd annual meeting of the association for computational linguistics: system demonstrations*, pages 55–60, 2014.
- [12] DG Maynard and Mark A Greenwood. Who cares about sarcastic tweets? investigating the impact of sarcasm on sentiment analysis. In *LREC 2014 Proceedings*. ELRA, 2014.
- [13] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*, 2013.

- [14] Rishabh Misra and Prahal Arora. Sarcasm detection using hybrid neural network. *arXiv preprint arXiv:1908.07414*, 2019.
- [15] Jeffrey Pennington, Richard Socher, and Christopher Manning. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 1532–1543, 2014.
- [16] Matthew E Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. Deep contextualized word representations. *arXiv preprint arXiv:1802.05365*, 2018.
- [17] Soujanya Poria, Erik Cambria, Devamanyu Hazarika, and Prateek Vij. A deeper look into sarcastic tweets using deep convolutional neural networks. *arXiv preprint arXiv:1610.08815*, 2016.
- [18] Tomáš Ptáček, Ivan Habernal, and Jun Hong. Sarcasm detection on czech and english twitter. In *Proceedings of COLING 2014, the 25th International Conference on Computational Linguistics: Technical Papers*, pages 213–223, 2014.
- [19] Antonio Reyes, Paolo Rosso, and Tony Veale. A multidimensional approach for detecting irony in twitter. *Language resources and evaluation*, 47(1):239–268, 2013.
- [20] Ellen Riloff, Ashequl Qadir, Prafulla Surve, Lalindra De Silva, Nathan Gilbert, and Ruihong Huang. Sarcasm as contrast between a positive sentiment and negative situation. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 704–714, 2013.
- [21] Stephen E Robertson and K Sparck Jones. Relevance weighting of search terms. *Journal of the American Society for Information science*, 27(3):129–146, 1976.
- [22] Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D Manning, Andrew Y Ng, and Christopher Potts. Recursive deep models for semantic compositionality over a sentiment treebank. In *Proceedings of the 2013 conference on empirical methods in natural language processing*, pages 1631–1642, 2013.
- [23] Oren Tsur, Dmitry Davidov, and Ari Rappoport. Icwsm—a great catchy name: Semi-supervised recognition of sarcastic sentences in online product reviews. In *Fourth International AAAI Conference on Weblogs and Social Media*, 2010.
- [24] Zichao Yang, Diyi Yang, Chris Dyer, Xiaodong He, Alex Smola, and Eduard Hovy. Hierarchical attention networks for document classification. In *Proceedings of the 2016 conference of the North American chapter of the association for computational linguistics: human language technologies*, pages 1480–1489, 2016.
- [25] Meishan Zhang, Yue Zhang, and Guohong Fu. Tweet sarcasm detection using deep neural network. In *Proceedings of COLING 2016, The 26th International Conference on Computational Linguistics: Technical Papers*, pages 2449–2460, 2016.

- [26] Xiang Zhang, Junbo Zhao, and Yann LeCun. Character-level convolutional networks for text classification. In *Advances in neural information processing systems*, pages 649–657, 2015.