# Fake News Pattern Recognition using Linguistic Analysis

BRAC
UNIVERSITY

Inspiring Excellence

**SUBMISSION DATE: 25.03.18**

**SUBMITTED BY:**

Amitabha Dey
Rafsan Zani Rafi
Shahriar Hasan Parash
Sauvik Kundu Arko

Department of Computer Science and Engineering

## Supervisor:

**Amitabha Chakrabarty, Ph.D**
Assistant Professor
Department of Computer Science and Engineering

# Declaration

We, hereby declare that this thesis is based on results we have found ourselves. Materials of work from researchers conducted by others are mentioned in references.

**Signature of Supervisor**                    **Signature of Authors**

_____

**Amitabha Chakrabarty, Ph.D**
Assistant Professor
Department of Computer Science and
Engineering
BRAC University

_____

**Amitabha Dey(13201070)**

_____

**Rafsan Zani Rafi(13201030)**

_____

**Shahriar Hasan
Parash(14101233)**

_____

**Sauvik Kundu
Arko(1320100055)**

# Acknowledgement

# ABSTRACT

The upsurge of fake news in social media calls attention to the erosion of long-standing institutional defense against misinformation in this Digital Age. In the wake of the 2016 Presidential Election in US, where social media played a crucial role in swinging votes, fake news has been a subject of increased discussion and debate. Social media used for news consumption has both its perks and disadvantages. In one hand it is relatively inexpensive and can be easily accessed but at the same time the likelihood of falling prey to fake news cannot also be disregarded. In this paper, we initially examine some of the existing technologies and frameworks that have been adopted to augment humans to make better decisions when it comes to recognizing news deception. We perform veracity assessment; conduct a comprehensive linguistic analysis on tweets to extract bag-of-words to perform Classification, specially structured around a specific target, in an attempt to find noticeable pattern in reliable and untrustworthy news. We discuss several surveys that have been undertaken in the past to help us present a comprehensive review of detecting fake news on social media. We later turn to the discussion of interconnected research domains and future research directions for constructing an ideal model for fake news detection system around social media. Although designing a fake news detector is not a straightforward problem, we propose a head-down operational guideline for a feasible fake news detecting system from a linguistic perspective.

**Keywords** - Deception, detection, social media, news verification, Bag of Words, linguistic analysis, semantic analysis, recognition, pattern, graph.

# TABLE OF CONTENTS

III

# LIST OF FIGURES

V

# LIST OF TABLES

# CHAPTER 1

# Introduction

In general humans are fairly ineffective at recognizing deception at least we have to take that as consideration [1]. Initially, most people are biased towards truth [2]; they tend to believe that the information they perceive is a fact and reliable. After that exhibition of general gullibility is also present among some people [3] and are strikingly acceptant to the concept those are not trivially perceivable. Advancement of technology made fake news from a minor internet sideshow to major electoral threat in a far greater speed than wildfire. Numerous approaches have been performed throughout the decade in search of obtaining a way to ensure news legitimacy. Immoral news presentation causes an immense of confusion in society. Even some of the fake news even triggered unpleasant incidents and cause losses of a huge amount of resources in terms of time and money. News represents a vital source of information along with knowledge for people. Eventually, in an era of world shaped by social media, the escalation of fake or hoax based news certainly has put traditional journalism sources and media system to a challenge. Initially, one of the challenges people face in detecting misinformation is that there does not yet exist a definite definition to explain fake news and the characteristics needed to determine articles legitimacy [4]. The task of detecting news to be faked is interpreted as the forecasting of the possibilities of being deceptive on purpose for a particular news article (news report, editorial, expose, etc.) [5]. Fake news is in a sense also entitled as fictional stories with a motive to deceive. By extension, fake news detection is the task of calculating the chances of an article being deceiving

or at least find the quantum of fake information represented [6]. Moreover posted articles or piece of information occurred in social media and political discourse consist a considerable power in forming people beliefs and opinions in distinct ways. As a result, their transparency with news is often compromised to amplify impact on society [7]. The culture of accepting news from social media is discernible. As an example, almost 62% of U.S. adults rely for news upon various social media based platform in 2016, while during 2012; only 49% of the adult populations affirm perceiving news from social media. Another report says that social media at this point of time out sell the very television itself as the considerable news source [5]. Results suggest in detecting lies in a text than chance there is a 4% margin of betterment, which was based on more than 200 experiments after complete meta-analysis [8]. In terms of helping to reduce the negative effects caused by fake news and also to welfare the public along with the news ecosystem - It's crucial that we should develop systems to detect fake news on social media or other online news portal automatically. In the latest time, diminishing trust in the mainstream media has been a pattern among people [1]. According to Gallup polls, only 40% population of adult Americans still have trust on their mass media sources to report the news fully, accurately and fairly [9].

## 1.1 Motivation

Fake news is defined as deliberately false information spread via print, broadcast, or online social media. Aside from reporting errors, in general, fake news does not come from established news sources. It is written with the intent to mislead in order to gain financially or politically. It is factually

incorrect and usually has sensational with headlines designed to grab attention. The shocking headlines and emotion-invoking text are deliberately composed to generate broad popular appeal and to encourage widespread sharing. One of the problems with discussing fake news is that it appears in multiple forms. It is authored for a variety of motives by an assortment of diverse individuals. To further complicate the definition, the term fake news has been usurped and is sometimes improperly used by individuals to disagree with or choose not to recognize the facts of a news story. So our motivation for this research is to understand the pattern of fake news and study it's characteristics by linguistic approach.

## 1.2   Thesis contribution

The objective of this research is to analyze the fake news, as the existing systems are mostly dependent on stance detection regarding this issue. Our research incorporated 'Bag of Words' to find out unique words and occurrence of those words in an article, which is well renowned technique for this type linguistic based analysis. Then we applied KNN algorithm over it for further analysis. In this research our approach was successful in order to differentiate legitimate tweets on 'Hillary Clinton' during US presidential election.

## 1.3   Thesis Outline

In the **Chapter 1** we basically discussed about our thesis motivation, methodology and introduction about our thesis topic.

In the **Chapter 2** we emphasized upon literature review and philosophical value of news. We also discussed about linguistics in this chapter.

In the **Chapter 3** our whole pattern analysis methods and findings are stitched together. This is the chapter which also gives us a whole understanding of our analysis method.

Lastly, in the **Chapter 4** our report concludes and we have also stated our future work plan here.

## 1.4 Methodology

In this portion we acknowledge stance detection though we did not incorporated it in our thesis and we discussed about KNN algorithm and Bag of Words.

### 1.4.1 Stance Detection

Stance detection comprises the estimation of the relative perspectives of two different text pieces on the same topic as described by [10]. Specifically, the task is to estimate the stance of a news headline, relative to the contents of a news article which can but does not have to address the same topic. Thus, the relative stance of each headline-article pair has to be classified as either unrelated, discuss, agree or disagree [11].

### 1.4.2 KNN Algorithm

KNN calculation is one of the least complex arrangement calculations and it is a standout amongst the most utilized learning calculations. KNN is a non-parametric, linguistic learning calculation. Its motivation is to utilize a database in which the information focuses are isolated into a few classes to anticipate the order of another example point. KNN has no model other than putting away the whole dataset, so there is no learning required. Efficient usage can store the information utilizing complex information structures like

k-d trees to influence hope to up and coordinating of new examples amid expectation productive.

### 1.4.3 Bag of Words

The Bag of Words is a rearranging portrayal utilized as a part of normal dialect preparing and data recovery (IR). Otherwise called the vector spaces demonstrate. In this model, content, (for example, a sentence or a record) is spoken to as the sack (multi set) of its words, slighting language structure and even word arrange yet keeping variety. The Bag of Words has been additionally utilized for linguistic analysis of different kinds. Bag of Words has pretty much helpful functionality for extracting any common pattern over words within a sentence.

# CHAPTER 2

# Literature Review

## 2.1  Philosophical Position

A spectacular form of understanding effective arguments is to the Aristotelian concepts of Ethos, Pathos, and Logos [12]. This demands a standard working rhetoric sense. The vitality of persuasive and legitimate writing is the capability to dissect and validate, or debunk the verbal aspects of other arguments.



**Figure. 2.1. Aristotle divided the aspects of persuasion into three categories: ethos (credibility), pathos (emotion) and logos (logic).**

As credibility refers to people believing who they trust, emotion and logic indicate a person's emotional connection and means of reasoning to convince one of a particular argument and/or speech.

As shown in Figure. 2.1,"Ethos" points to the writer's "ethics," which represents a writers control or character to entente with a topic. Ethos is basically stands for the credibility measure of the deliverer. For the love of engaging audience to a certain topic, the character presenting the information must be established thus considered to be trusted, also an experienced one also. This is basically ethical quantifier. Pathos points to the arguments emotional plea or the writings ability to establish an ingrained relation with the receiver. Although it does not have that much effect to shift the readers emotional spectrum. It is more about moving their mind. This creates an idea of emotion being constructed by ethos with motion. Several times this appeal is how a writer will make an argument important to a reader. Logos associated with arguments logical standing. An effective argument will acknowledge objectivity and other supporting details to prove authors claims or standing. A testimony from superiors with cautiousness of the writer to choose better evidence in order to back up personal claim is a vital part of it. In most cases a valid writing is prone to be well organized and written with sheer skill.

## 2.2 Linguistics

Languages represent sets of signs. Signs fuse an exponent (a sequence of letters or sounds) with a deep meaning [13]. Grammars are system to generate signs from more basic signs. Signs combine a form and a meaning, and they are identical with neither their exponent nor with their meaning. Language is a means to connect with people, it is a semiotic system. By that we basically mean that it is a set of signs. Its A sign is a pair consisting—in

the words of Ferdinand de Saussure—of a signifier and a signified. We like to call the signifier the exponent and the signified the meaning. In linguistics, language signs are constituted of four different levels, not just two: phonology, morphology, syntax and semantics. Semantics deals with the meanings (what is signified), while the other three are all concerned with the exponent. At the lowest level we find that all is composed from a small set of sounds, or—when we write—of letters. The minimal parts of speech that bear meaning are called 'morphemes'. Often, it is tacitly assumed that a morpheme is a part of a word; bigger portions are called idioms. Word has meaning, its sound structure, its morphological structure and its syntactic structure. The levels of manifestation are also called strata [14]. There are levels of linguistic analysis which are;

1. **Phonetics:** Phonetics is the study of production, transmission and perception of speech sound. It is concerned with the sounds of languages, how these sounds are articulated and how the hearer perceives them.

2. **Morphology:** It is study of word formation and structure. It studies how words are put together from their smaller parts and the rules governing this process.

3. **Lexicology:** It is study of words. We study word-formation and world classes. Lexeme is the smallest unit of Lexis.

4. **Syntax:** It is the study of sentence structure. It attempts to describe what grammatical rules is in particular language.

5. **Semantics:** It is the study of meaning in language. It is concerned with describing how we represent the meaning of word in our mind how

we use this representation in constructing sentence. It is based largely on the study logic in philosophy [15].

It also associates insights into the nature of language variation (i.e. dialects), language shifting over period, how language has been processed and saved in the brain, and in which way it is acquired by early age population. All of these data were studied and analyzed for over a decade by the University of Arizona's Department of Linguistics. Although linguistics is still vaguely unfamiliar to the educated portion of the population, it is an expanding and enthusiastic field with a swelling major impact on other fields as diverse as psychology, philosophy, education, language teaching, sociology, anthropology, computer science, and artificial intelligence. A researcher with an interest in linguistics can select among a diverse range of study paths. Some of these are renowned for various important research modification purposes. A point to be noted is, different study paths will be fruitful from different course concentrations, so it's a generous idea to explore these zones of linguistics which can assist to learn about the patterns those fake news follow and identify them before they can cause any negative impact on the whole society and the social media by which it is stitched in the modern civilization. Linguistics based approach were introduced a long ago in the sector of machine learning and it is thriving since that time.

# CHAPTER 3

# Fake News Pattern Analysis & Result

## 3.1 Veracity Assessment

This portion is about trust ratio and list of credible sources we derived from our study.

## 3.1.1 Most to Least Trusted News Outlets in USA

We primarily examined two comprehensive surveys in order to establish the credibility of US News sources. The first study was initiated by pew, a well established Research Center in 2014. Pew surveyed a representative sample of randomly selected Americans, polling nearly 3,000 people in 2014. The second study was an online survey in 2017, initiated by 28 different news providers in US. It attracted almost 9,000 respondents, who were asked to name three news sources they trusted and three they did not [16]. These two major surveys have been widely acknowledged by several independent professional journalist societies.

**Figure. 3.1.1(a) To generate these graphs, we used matplotlib and xlrd Python libraries to create a system that can fetch data from a dynamic excel file and display bar charts respectively. The Economist had the highest trust rating overall, while Occupy Democrats had the least.**

With the aid of these surveys, we constructed our own trust-ratio scale of the most and least credible news outlet in America. Graphs representing our trust-ratio scale are shown in Figure. 3.1.1(a) and (b).



**Figure. 3.1.1(b) News sources vs Trust Ratio graph, generated from the table in Table 3.1.1.**

A room for other news sources to be incorporated in the trust ratio is always open and appreciated in the system. Some of the well recognized ones are pre loaded as starter for the system due to comprehensive standings of reliable news sources is different for several studies. So eventually we construct an average scale of reliable sources stated in all studies we were provided with.

| Legends | Source | Score |
|---|---|---|
| N1 | The Economist | 10 |
| N2 | Reuters | 10 |
| N3 | BBC | 10 |
| N4 | Guardian | 10 |
| N5 | WSJ | 10 |
| N6 | PBS | 9 |
| N7 | NPR | 9 |
| N8 | ABC News | 9 |
| N9 | Politico | 9 |
| N10 | Associated Press | 8 |
| N11 | NBC | 8 |
| N12 | CNN | 8 |
| N13 | USA Today | 8 |
| N14 | LA Times | 8 |
| N15 | Google News | 8 |
| N16 | The NYT | 7 |
| N17 | Denver Post | 7 |
| N18 | Washington Post | 7 |
| N19 | Atlantic | 5 |
| N20 | CBS | 5 |

**Table 3.1.1. News sources with their respective Trust Ratio Score.**

Initial finding of this research was to filter top 20 'United States of America' credible sources prepared by studying survey and research report done by transparent criticism. Later integration of survey results assist to structure a trust ratio scale, which is illustrated in figure 3.1.1(b).

**3.1.2 Determining Alexa Rank Score via domain parsing**

Alexa Rank, designed by Amazon, is a metric that ranks website in order of popularity in terms of web traffic. It takes into account an estimation of the average daily unique visitors and the underlying algorithm redresses for

various possible biases. For this, we designed a module using urllib parser to extract the domain name from the link of an article and return its respective Alexa rank. An example of what this module accomplishes is shown in Figure. 3.1.2. The Alexa rank is demographically depended and gives different values from region to region. Due to this credibility of this score is minor weighted and treated as optional.

```
Enter URL : http://www.bbc.co.uk/news/world-europe-43389407
http://www.bbc.co.uk/
Alexa rank :
101
```

**Figure. 3.1.2. Output for a single source of this module.**

### 3.1.3 Extraction of Hash tags

Hash tags have become a popular way of tagging content on Social Media, and attaching hash tags to a piece of content can dramatically increase its visibility on social media. One of our modules returns all the hash tags associated with a particular article from analyzing the entire webpage and creates a word-cloud. An example of what this module accomplishes is shown in Figure. 3.1.3. Hash tags can be used as a key element to fetch important information about any particular topics. Though hash tags can be some time confusing and may not be efficient for information search but apart from that hash tags are ingenious way to extract co related key words from a piece of text. These key words can act as summary and notable nouns or actions present in the news.

**Figure. 3.1.3.  Hash tag word cloud generation by world of hash tags.**

Key points from a post are accumulated in the word cloud which is derived from hash tags. This technique indeed assists the perceiver to learn about the acting characters, related places, performed actions and other trigger words mentioned within a post or article.

## 3.2 Headline Analysis

In this portion we discussed about headline analysis and other methods related to that.

### 3.2.1 POS Tagging Using NLTK

The NLTK framework is used for processing natural languages and providing comprehensive support for various NLP related tasks [14]. We

used NLTK to tag parts-of-speech to each token in an attempt to recognize words of importance. The output generated by module is passed down to the next modules. An example of what this module accomplishes is shown in Figure. 3.2.1. These extraction is necessary to pass them in the event registry module and later it can be also be analyzed with association of dandelion API which is mainly used for text similarity analysis. So undeniably NLTK provides the skeleton of our system to drive forward. NLTK pluck out parts of speech and assists to structure the skeleton of the whole system. All  of the words extracted as parts of speech is analyzed with Bag of Words which provides another point of view to the analyzing method and some words such as preposition, articles, punctuations are also filtered out along the process .

```
[('The', 'DT'), ('United', 'NNP'), ('States', 'NNPS'), ('could', 'MD'), ('have',
    DT'), ('permanent', 'JJ'), ('division', 'NN'), ('of', 'IN'), ('Europe', 'NNP')
     'VBN'), ('complicit', 'NNS'), ('in', 'IN'), ('the', 'DT'), ('oppression', 'NN
    ('.', '.')]
[('Today', 'NN'), (',', ','), ('having', 'VBG'), ('come', 'VBN'), ('far', 'RB'), (
     'JJ'), ('historical', 'JJ'), ('journey', 'NN'), (',', ','), ('we', 'PRP'), ('
     ':'), ('Will', 'MD'), ('we', 'PRP'), ('turn', 'VB'), ('back', 'RP'), (',', ',
    'well', 'RB'), ('?', '.')]
[('Before', 'IN'), ('history', 'NN'), ('is', 'VBZ'), ('written', 'VBN'), ('down',
    , (',', ','), ('it', 'PRP'), ('is', 'VBZ'), ('written', 'VBN'), ('in', 'IN'),
[('Like', 'IN'), ('Americans', 'NNPS'), ('before', 'IN'), ('us', 'PRP'), (',', ','
    show', 'VB'), ('that', 'DT'), ('courage', 'NN'), ('and', 'CC'), ('we', 'PRP'),
    well', 'RB'), ('.', '.')]
[('We', 'PRP'), ('will', 'MD'), ('lead', 'VB'), ('freedom', 'NN'), ("'s", 'POS'),
[('We', 'PRP'), ('will', 'MD'), ('compete', 'VB'), ('and', 'CC'), ('excel', 'VB'),
    global', 'JJ'), ('economy', 'NN'), ('.', '.')]
[('We', 'PRP'), ('will', 'MD'), ('renew', 'VB'), ('the', 'DT'), ('defining', 'VBG'
    NNS'), ('of', 'IN'), ('this', 'DT'), ('land', 'NN'), ('.', '.')]
[('And', 'CC'), ('so', 'RB'), ('we', 'PRP'), ('move', 'VBP'), ('forward', 'RB'), (
    about', 'IN'), ('our', 'PRP$'), ('country', 'NN'), (',', ','), ('faithful', 'J
```

**Figure. 3.2.1. Our module performs POS tag on every token.**

### 3.2.2 Named Entity Recognition Using NLTK

This is the first step after POS tagging towards information extraction from unstructured data. Named Entity Recognition is part of the information extraction. It is known as entity identification and entity extraction [5]. Named Entity Recognition can be used to automatically populate a legal ontology from legal texts following ontology learning [17]. [18] Shows how web resources such as Wikipedia and Wiktionary can be used in combination with a domain corpus, a general purpose named entity tagger and a seed or base ontology to derive domain ontology. We used NLTK to construct a NER tree. An example of what this module accomplishes is shown in Figure. 3.2.2.



**Figure. 3.2.2. This module extracted real-world entity from the headline (Person, Organization, etc)**

### 3.2.3 Searching Global News Using EventRegistry

Event Registry is a system that can analyze news articles dropped in to it and search for world news events mentioned in them. The system is capable to identify a bulk of articles that has coherence with the same event. It is able

to identify groups of articles in different languages that describe the same event and represent them as a unique event. From articles in each event it can then extract events significant information, such as event related place, time and date, personalities who are related and what is it talking about. Extracted information is kept in a database. A dedicated interface for the users is available that allows users to search for events using advanced search options, to visualize and aggregate the search results, to analyze individual events and to identify correlated events [19]. By using Named Entity Recognition to collect searchable and meaningful token from the headline of the article and parallelly performed queries in Event Registry based on those extracted tokens to fetch news published by various news outlets or portals around the globe on that particular topic. An example of what this module accomplishes is shown in Figure. 3.2.3.

```
q = QueryArticlesIter(keywords=QueryItems.OR(["Hillary", "Russia"]))
```

```
{'id': '11779704', 'uri': '833562650', 'lang': 'eng', 'isDuplicate': True, 'date': '2018-03-17', 'time': '19:53:00', '
    dateTime': '2018-03-17T19:53:00Z', 'sim': 0, 'url': 'http://www.vvdailypress.com/zz/news/20180317/
    latest-ex-cia-director-accuses-trump-of-moral-turpitude?rssfeed=true', 'title': "Trump lauds firing of ex-top FBI
    official as 'great day'", 'body': 'Attorney General Jeff Sessions fired Andrew McCabe, a former FBI deputy
    director long scorned by Trump\n\nWASHINGTON -- In what President Donald Trump called "a great day for Democracy,"
    Attorney General Jeff Sessions fired Andrew McCabe, a former FBI deputy director long scorned by Trump, two days
    before McCabe\'s scheduled retirement date, acting on the recommendation of bureau disciplinary officials.\n\n
    McCabe suggested the move was part of the Trump administration\'s "war on the FBI." Trump tweeted in praise of
```

**Figure. 3.2.3 Based on the named entity recognitions we extracted in our previous module, we used the noun-phrases as query. A query of "Hillary" and "Russia" as such, for example, returned all the related news covering the topic from around the world.**

### 3.2.4 Comparing Text Similarity

The Dandelion API consists of a domain offering endpoints, for distinct text analysis tasks. In specific, this API offers semantic analysis features for: Entity Extraction, Text Similarity, Text Classification, Language Detection, and Sentiment Analysis [20]. Dandelion offers an extent of knowledge graph of locations, events, organizations, persons and other information.



**Figure. 3.2.4. Entity extraction using Dandelion.**

In common data extracted from various different data origins are sewed altogether in a single substantial knowledge graph by Dandelion, and it also provides a set of APIs over it [21]. Dandelion API is a structure crafted to direct this programmability provocation for data-parallel applications. Dandelion gives an insight to a unified programming model for

heterogeneous systems that extent a divergent array of execution contexts including CPUs, GPUs, FPGAs, and the cloud. According to L.Ross et. al. Dandelion implements a single machine abstraction: the programmer composes continuous code in a high-level programming language [22]. We used Dandelion to compare the Syntax and Semantic similarity between two excerpts.

This actually shows the inter connectivity of news to determine the legitimacy of it. In most cases Dandelion takes a piece of text and compares its structure with another. This phenomenon also works well for news comparison if key characters of the news can be passed in this particular API.

```
Please enter Text 1 :Hillary Clinton warns Britain on potential trade deal with Trump
Please enter Text 2 :Hillary Clinton says President Trump is continuing to "ignore and surrender" to Russian
meddling, asking whether the Republican who defeated her in the 2016 election will
 protect the country
Similarilty : 71.39999999999999%
```

**Figure. 3.2.4(a). This module can compare the syntax and semantic similarity between two tweets. The level of accuracy is not necessary consistent, but this is a marker that needs to be kept track of.**

### 3.2.5 Performing Image Search

This module collects all the noun phrases extracted from the headline of the article by TextBlob, which is a Python library for processing textual data, and we parse and concatenate to build an URL to display all the Google Images based on the query. An example of what this module accomplishes is

shown in Figure. 3.2.5. This will also act as a visual aid to what sort of content the user is consuming.



**Figure. 3.2.5. Based on the headline, this module is able to fetch relevant images from Google.**

## 3.3 Linguistic Approach

According to Roxana et. al. [23] Languages represent varieties domains of signs. Signs include an exponent (a continuous sequence of letters) with an implication. Grammars open up ways to craft signs from more general signs. Signs fuse a form and a content, and they are identical with neither their exponent nor with their definition. Most of false claims show usage of language in such a strategic way that might help to avoid being unmasked as fake. In addition to the action to control their speech, language leakage occurs with particular verbal aspects that are tougher to watch over such as frequencies and patterns of pronoun, conjunction, and use of negative

emotional texts [24]. Dramatization of a news story may vaguely depend on usage of Subjective words [15]. Again, in the domain deception detection research satire is famous for being an attractive subject: it is represents a twisted type of deception that is in a way more intentional and incorporates cues revealing its deceptiveness of its own [8].

### 3.3.1 Data Representation

Natural language processing (NLP) systems attain strings of words (sentences) as their input and generate a structured representations pull out the meaning of those strings as their outcome. The goal of linguistic approach is to look for language leakage or so called Predictive deception cues found in the content of a message [25].Perhaps the simplest method of representing texts is the bag of words approach, which regards each word as a single, even significant unit. In the bag of words approach, individual words on n-grams (multi word) frequencies are aggregated and analyzed to reveal cues of deception [26]. Respective lexical cues are tagged to words as a function of the whole process e.g. parts of speech [27], affective dimensions [28] or location-based words [29] are recognized as ways to reveal linguistic indication of deception by supplying subsets of frequency.
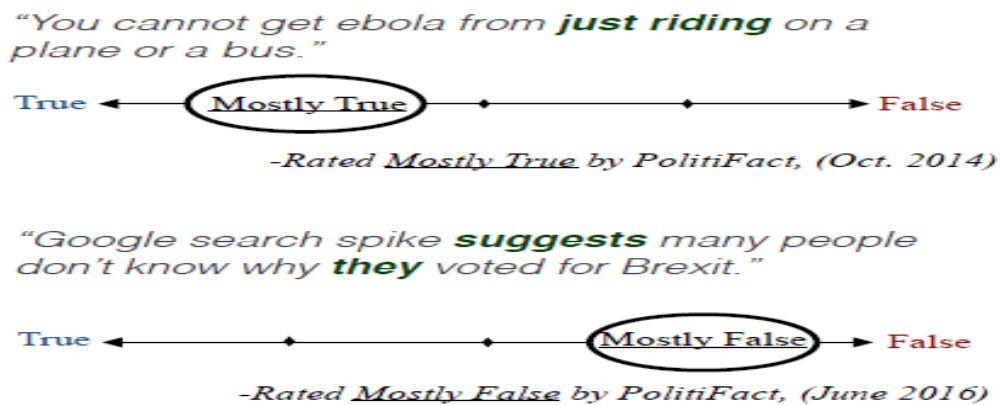
### 3.3.2 Deep Syntax

The task of news analyzing is not always well enough in forecasting deceptive information. For purpose of speculating norms of deception deeper language structures (syntax) have been closely examined [30]. Outside assessment tools e.g. the Stanford Parser, [10; 31], AutoSlog-TS syntax analyzer [15] and others aid in the process constructing automatic

deep syntax analysis. Solely, syntax analysis might not have adequate ability of recognizing deception, and researches generally merge this viewpoint along with other linguistic or network resulation methods [32].

### 3.3.3 Ambiguity in Language

Deception Detection Psycholinguistic work in interpersonal deception theory [33] has hypothesized that a particular pattern of writing can suggest the signs of a writer willing to purposefully obscure the fact, as shown in Figure.11.



**Figure. 3.3.3. Based on the headline, this module is able to fetch relevant images from Google.**

Hedge words and other vague qualiers [34; 23], e.g. this may increase indirectness to a position in the writing that bemuse its meaning. Again satirical statements may trigger more generalized traits that reveal its deprivation from the truth because the objectivity of sarcasm itself is to be amusing and understandable at the same time towards at least a few set of readers to recognize the humor of it [24]. Again articles derived from The Onion and other satirical news sources are often restated in other news

outlets or shared on online platforms as if the stories were factual [7]. Traditionally, satire has been distributed into two styles: Juvenalian, the more tending to hostility between two and Horatian, the more lively [24]. Juvenalian type of satire is specified by repulsion and sarcasm where Horatian satire by contrast on nature tends towards teasing and mockery [33].

## 3.4 Tweet Analysis

This portion is about tweet analysis and our operations on a tweet.

### 3.4.1 Stance Detection for Twitter

Stance detection for tweets involves detecting if the tweet is in FAVOR or AGAINST a particular target which can be a person, a trending topic, etc [1]. The idea is to collect two sets of dataset on a particular target: One should be labeled as "Legitimate" and the other as "Fake". For our research, we decided on our target as "Hillary Clinton", since it was one of those topics that have been intensely spoken about in the last year. We collected 80 tweets on "Hillary Clinton" tweeted by the five most credible news sources listed in our aforementioned report. These sources are most likely to hold credibility in terms of representing news as it is without any kind of biasness. Later from those collected tweets we can also run sentiment analysis module to extract polarity from given tweets and subjectivity of those tweets also. A sample of the dataset is shown in Table 3.4.1.

| Date Added | Source | Tweets |
|---|---|---|
| 14 March, 2018 | CNN | President Trump slams Hillary Clinton during a speech in California: "You wouldn't be going to Mars if my opponent won" |
| 12 March, 2018 | Reuters | Hillary Clinton sees the sights while in India for a private visit.<br>More video: |
| 22 February, 2018 | The Economist | One former troll described how they hired a Hillary Clinton lookalike prostitute to have sex with a black man on camera |
| 19 February 2018 | The New York Times | Using American bank accounts, drivers' licenses and disposable phones, about 80 Russians worked to disparage Hillary Clinton, promote Trump and sow political discord |
| 11 November, 2017 | The Guardian | Woman jailed over Hillary Clinton documentary fraud |

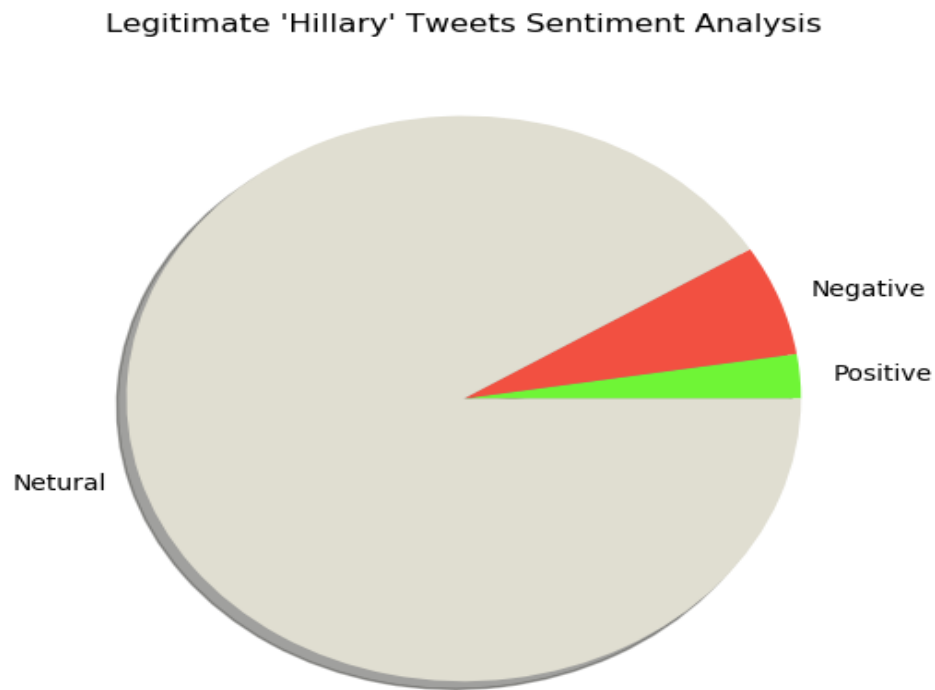**Table  3.4.1. Sample of tweets collected from news sources with high trust ratio – Labeled "Legitimate".**

## 3.4.2 Sentiment Analysis

We then performed Sentiment analysis on all the 80 tweets, measuring Polarity (-1 to 1) and Subjectivity (0.0 to 1.0, where 0.0 is very objective and 1.0 is very subjective.).

| Tweets | Polarity | Subjectivity |
|---|---|---|
| Hillary Clinton told a receptive audience in India that while she thought President Donald Trump played to some of Americans' worst fears, he does not reflect the country as a whole  -- and said the US "did not deserve" his presidency | -0.4 | 0.7 |
| Hillary Clinton says President Trump is continuing to "ignore and surrender" to Russian meddling, asking whether the Republican who defeated her in the 2016 election will "protect the country" | 0.0 | 0.0 |
| Hillary Clinton attacks proposed Trump budget cuts as 'cruelty' | 0.0 | 0.0 |
| In October 2017 Hillary Clinton told us that Democrats should stand up for what they believe in, even if it costs them votes | 0.0 | 0.0 |
| Congressional Republicans are launching joint investigations into multiple issues concerning Hillary Clinton | 0.0 | 0.0 |

**Table  3.4.2. "Hillary" tweets with its polarity and subjectivity.**

In Table 3.4.2, we can observe that the "Legitimate" tweets on "Hillary" are mostly Neutral in terms of Polarity. We can therefore hold it as one our markers for constructing predictive model. The pie chart in Figure 3.4.2(a) tells us that the "Legitimate" tweets on "Hillary" are mostly Neutral in terms of Polarity. We can therefore hold it as one our markers for constructing predictive model.



**Figure. 3.4.2(a). Sentiment analysis of Legitimate "Hillary" tweets at a glance.**

President Donald Trump is infamous for his highly polarized tweets worldwide. Most of Trump's tweets exhibit a lot of linguistic characteristics

usually observed in otherwise deceptive news. His tweets are ambiguous in nature, abusive and mostly subjective. We collected 50 tweets of Donald Trump on "Hillary" from the same time period of when we collected the "Legitimate" tweets. We performed the above sentiment analysis on all of these tweets as well and noted them down accordingly. The characteristics of Trump's tweets can be reflected in the pie chart in 3.4.2(b), generated by analyzing his tweets. Donald Trump's hoax based tweets seem to provide negative polarity more than the legitimate Hillary tweet dataset.

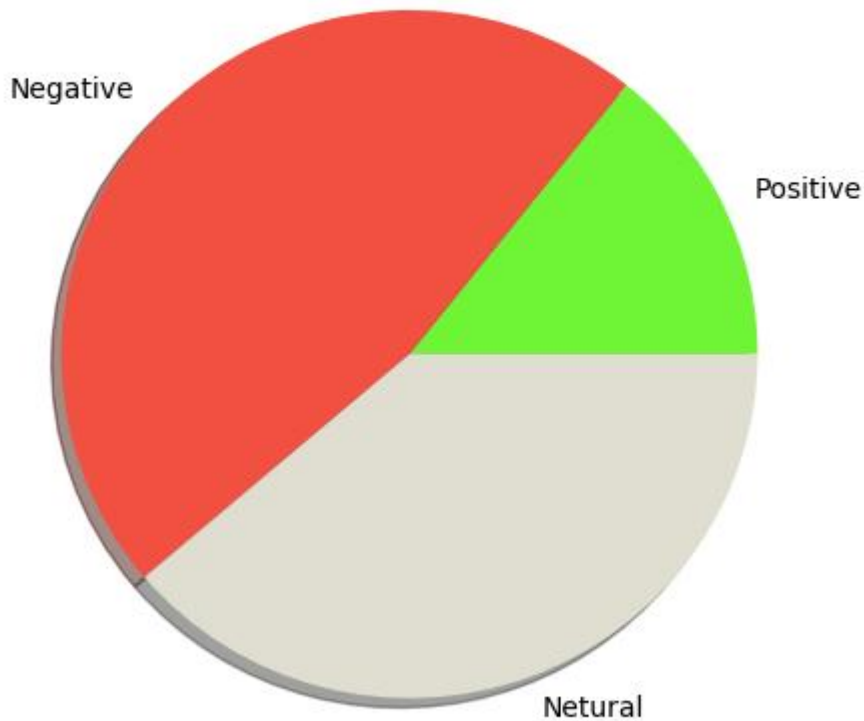Donald Trump's polarized 'Hillary' Tweets Sentiment Analysis



**Figure.  3.4.2(b). Sentiment analysis of Donald Trump's polarized "Hillary" tweets.**

### 3.4.3 Result & Findings

Perhaps the simplest method of representing texts is the bag of words approach, which regards each word as a single, but significant unit. We created a BoW model for the two datasets that we prepared in an attempt to find which word frequencies in each categorized tweets. We filtered through the 50"Hillary" tweets made by Donald Trump and appended only the noun-phrases in our corpus. We then created the BoW model based on the corpus. Our results are shown in figures.



**Figure. 3.4.3 (a). Every index represents an unique word from the entire corpus of tweets. 0 means there are no occurrence of the word, 1 means the word occurred once, and so on.**

Above figure is the illustration of BoW model matrix which stands for words occurrence in a certain tweet. Though frequency is not measure in this particular part but unique word occurrence is the main finding here. Triggered words can be shown to be present in a tweet by this array (first unique word index is generated) then their unique occurrence in every tweet for total dataset can also be derived by this. This is mandatory for pattern analysis in the further parts.

```
{'crooked': 17, 'hillary': 36, 'clinton': 10, 'was': 88, 'john': 40, 'podesta':
    62, 'big': 6, 'money': 51, 'russia': 71, 'dollars': 26, 'wife': 90, '
   political': 63, 'run': 70, 'drain': 29, 'swamp': 80, 'clintons': 11, 'dems'
   : 21, 'ads': 0, 'facebook': 32, 'media': 50, 'they': 83, 'wow': 91, 'james'
   : 39, 'comey': 13, 'terrible': 81, 'candidate': 9, 'hits': 37, 'dnc': 25, '
   fbi': 34, 'director': 24, 'free': 35, 'pass': 57, 'bad': 3, 'deeds': 19, '
   phony': 60, 'new': 53, 'polls': 64, 'fake': 33, 'intelligence': 38, 'bill':
    7, 'uranium': 87, 'russian': 72, 'speech': 76, 'korea': 43, 'nukes': 54, '
   election': 30, 'loss': 46, 'steve': 78, 'bannon': 4, 'thanks': 82, 'donald'
   : 27, 'authorities': 2, 'whereas': 89, 'special': 75, 'council': 15, '
   pocahontas': 61, 'primaries': 66, 'lets': 45, 'dept': 23, 'confirms': 14, '
   many': 49, 'real': 69, 'story': 79, 'collusion': 12, 'donna': 28, 'book': 8
   , 'primary': 67, 'crazy': 16, 'bernie': 5, 'justice': 41, 'department': 22,
    'sorry': 74, 'paul': 59, 'manafort': 48, 'trump': 85, 'presidential': 65,
   'movie': 52, 'star': 77, 'and': 1, 'kasich': 42, 'passion': 58, 'democratic
   ': 20, 'party': 56, 'sanders': 73, 'unfair': 86, 'obama': 55, 'lynch': 47,
   'law': 44, 'enforcement': 31, 'decisions': 18, 'purposes': 68, 'totally':
   84}
```

**Figure. 3.4.3 (b). Vocabulary corpus for the 50 Trump's "Hillary" tweets. It represents all the unique noun-phrases present in these tweets with their associated unique id number.**

```
word  6  occurs  3  times
word  7  occurs  1  times
word  8  occurs  1  times
word  9  occurs  1  times
word 10  occurs  4  times
word 11  occurs  1  times
word 12  occurs  1  times
```

**Figure.  3.4.3 (c). This tells us the occurrence of each word in all of the 50 tweets.**

Throughout the above diagrams the process of collecting each unique (Noun) words occurrence are visualized. By studying unique words

recursive occurrence in both fake and legitimate news a pattern can be analyzed. This linguistic approach give an insight to a generalize framework followed by deceiving or fake news and construct a differentiation between real to fake. Through calculating word frequency it can be perceived that the usage of a certain type word occurrence in a text. For example a deceiving word like 'big', 'huge', 'crooked', 'bad' etc. is used to hide quantum of any quality. Most times any article or informative text containing these type words are meant to create deception and hide factual information. If frequent used words in hoax or fake news can be identified then this appearance count can assist the hoax detection. More appearance of course means more likelihood to be fake news. We performed above modules on a particular topic which was about "Hillary Clinton" and "Donald Trump" tweets. Other diverse topics can also be analyzed using the same systems; here "Hillary" tweets are considered to be legitimate and labeled real, where "Trump" is considered as biased and labeled fake. Labeled dataset is also a key factor here because it saves an enormous time to screen out distortions in the data sets and provides a kick start to machine learning process. Though for huge datasets it is not that much feasible and is time consuming. We were also able to find out deceptive words and their interpretation in the sentence which is role player in sentiment analysis along with subjectivity. Larger dataset analysis would help to accomplish a dictionary of deceptive words which's presence can help to find out credibility of a news.

Our findings from this analysis is that Trump's polarized tweets were very subjective in nature. He often used superlative words and the word "big" frequented in his tweets which is not specific in nature, to describe a

situation and "phony" and "crooked", which is abusive in nature, to describe Hillary Clinton. We performed the same comprehensive analysis on the 80 Legitimate"Hillary" tweets. Our results are shown in the figure below. Based on our findings, we can assert that legitimate tweets in general are more specific and more objective in nature. They are rather critical than abusive. They refrain from using strong words if not absolutely necessary and would only do to quote somebody.

```
{'trump': 218, 'hillary': 98, 'clinton': 38, 'california': 29, 'mars': 135, '
  mike': 140, 'conaway': 43, 'texas': 206, 'republican': 183, 'intelligence':
  107, 'committee': 41, 'russia': 185, 'democratic': 53, 'presidential': 166
, 'nominee': 148, 'receptive': 179, 'audience': 14, 'india': 106, 'donald':
  63, 'presidency': 164, 'http': 102, 'cia': 37, 'gop': 90, 'rep': 182, '
chris': 35, 'stewart': 202, 'community': 42, 'assessment': 13, 'russian':
186, 'president': 165, 'putin': 175, 'private': 169, 'visit': 230, '
internet': 108, 'trolls': 217, 'video': 228, 'game': 87, 'hilltendo': 99, '
election': 67, 'senior': 192, 'adviser': 5, 'hud': 103, 'false': 73, '
conspiracy': 45, 'theory': 209, 'campaign': 31, 'chairman': 34, 'satanic':
189, 'cnn': 39, 'kfile': 122, 'uranium': 225, 'congress': 44, 'bill': 21, '
foundation': 82, 'opinion': 154, 'try': 219, 'open': 153, 'arms': 12, '
janus': 113, 'davidrivkin': 49, 'andrewmgrossman': 11, 'lookalike': 130, '
prostitute': 173, 'black': 23, 'man': 133, 'american': 10, 'bank': 17, '
accounts': 1, 'drivers': 66, 'licenses': 127, 'disposable': 59, 'phones':
159, 'political': 161, 'discord': 58, 'action': 3, 'school': 190, 'shooting
': 197, 'parkland': 157, 'florida': 80, 'november': 149, 'vote': 231, '
defendants': 52, 'prison': 168, 'uniform': 223, 'read': 178, 'rallies': 177
, 'freelance': 84, 'journalist': 119, 'christopher': 36, 'steele': 201, '
justice': 121, 'department': 54, 'watchdog': 232, 'bipartisan': 22, 'group'
: 92, 'fbi': 74, 'actions': 4, 'emails': 69, 'sexual': 195, 'harassment':
95, 'surprise': 203, 'fire': 77, 'fury': 86, 'grammys': 91, 'unexpected':
222, 'cameo': 30, 'us': 226, 'ambassador': 9, 'nikki': 147, 'haley': 94, '
audio': 15, 'version': 227, 'manager': 134, 'york': 238, 'breaking': 25, '
newly': 143, 'text': 207, 'exchanges': 70, 'officials': 152, 'email': 68, '
investigation': 110, 'special': 198, 'prosecutor': 171, 'recovers': 180, '
messages': 139, 'll': 129, 'candidates': 32, 'texts': 208, 'high': 97, '
level': 126, 'crucial': 47, 'donors': 65, 'treatment': 216, 'george': 88, '
papadopoulos': 156, 'foreign': 81, 'policy': 160, 'australian': 16, 'aide':
  7, 'dirt': 57, 'accuses': 2, 'deputy': 55, 'director': 56, 'october': 150,
  'agent': 6, 'mary': 137, 'beard': 18, 'creep': 46, 'little': 128, 'rock':
184, 'jeff': 114, 'sessions': 193, 'federal': 75, 'prosecutors': 172, '
related': 181, 'joe': 116, 'biden': 20, 'woman': 235, 'documentary': 62, '
fraud': 83, 'pundits': 174, 'donna': 64, 'brazile': 24, 'satan': 188, '
```

**Figure. 3.4.3(d) Vocabulary corpus for the 80 Legitimate"Hillary" tweets. It represents all the unique noun-phrases present in these tweets with their associated unique id number.**

K-nearest neighbors (KNN) is a classification technique used for both classification and regression predictive problems. At its most basic level, it is essentially classification by finding the most similar data points in the

training data, and making an educated guess based on their classifications. The steps of KNN algorithm are:
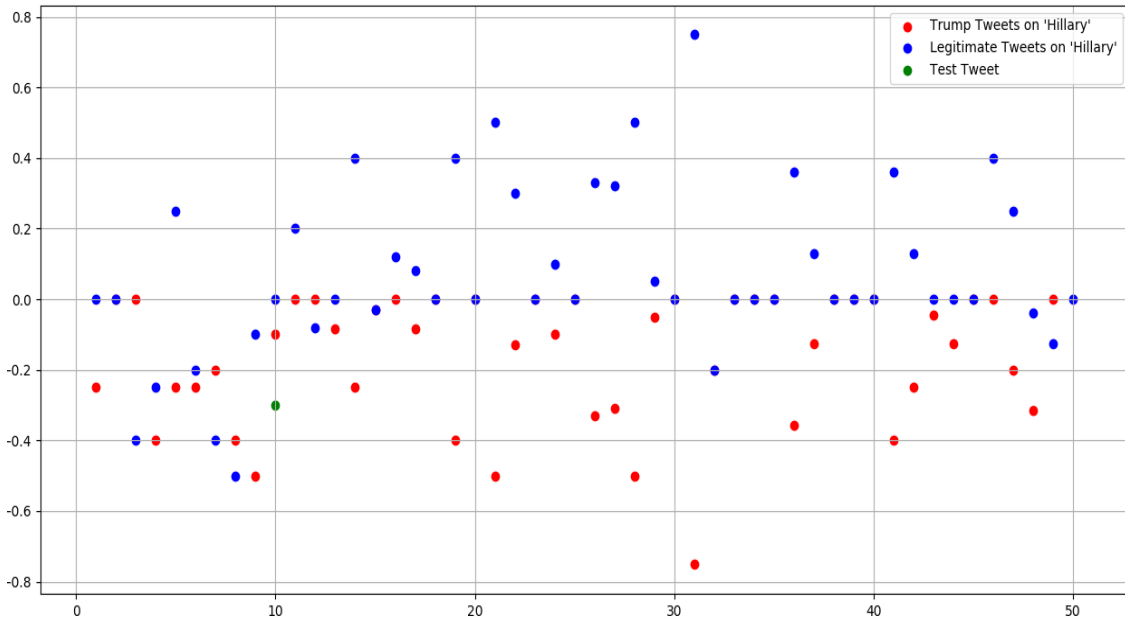
1.  Computing a distance value between the item to be classified and every item in the training data-set.
2.  Calculate the k closest data points (the items with the k lowest distances).
3.  Conducting a "majority vote" among those data points — the dominating classification in that is the final classification.

There are many different ways to compute distance, as it is a fairly ambiguous notion, and the proper metric to use is always going to be determined by the data-set and the classification task. Two popular ones, however, are Euclidean distance and Cosine similarity. We went for **Euclidean distance**.

$$E(x, y) = \sqrt{\sum_{i=0}^{n}(x_i - y_i)^2} \qquad \dots\dots\dots\dots (1)$$

Euclidean distance is essentially the magnitude of the vector obtained by subtracting the training data point from the point to be classified. In the figure below we plotted fake and legitimate tweets and show the standing of the test tweet. As it was labeled fake and machine was not trained for test tweet sets.

**Figure 3.4.3(e): Polarity vs. Tweets. 'Red' plots represent Trump's tweets and 'Blue' plots represent 'Legitimate' tweets. The single green plot represents our test data.**

The polarity of our randomly selected tweet on "Hillary" was -0.3. The application of KNN algorithm required us to first find the k nearest neighbors of our test data. We decided on k=3, i.e. the three nearest neighbors. By applying the Euclidian equation, we calculated the distances and located its three nearest neighbors. For this particular sample, out of the 3 nearest neighbors, 2 were labeled fake and 1 was labeled legitimate. The test sample tweet was 'Fake'. Our system gave a result for particular tweet where for known labeled tweet was not unveiled to machine but still it were able to find out that given news article pattern was most liked to be fake and voted 2/3 likelihood of being fake. The accuracy can be fine-tuned by increasing the number of tweets in our dataset.

# Chapter 4
# Conclusion

In this paper, we discussed the computational linguistics implementations we have used to perform linguistic analysis on tweets to observe patterns exhibited by legitimate and fake or ambiguous news. We have designed separate independent modules that are aimed to assist humans in making better decisions of detecting deception in news. We deconstructed the grammar of the tweets for in-depth analysis and constructed a comprehensive BoW model based on the categorized labeled tweets. We have compared how polarity and subjectivity varies between legitimate and polarized tweets, considering the topic of the tweets to be same. Our future research on this is going to be about performing in-depth stance detection analysis on top of the BoW models that we constructed in this research. Our framework can be used on other specified topic to create more learning assisted atmosphere for machine itself and also other diverse patterns can be learn by it and implemented for more efficient prediction in upcoming future days. Introduction of satire based linguistic to machine will be a big challenge to our future expansion of this model but still it can make our framework capable of analyzing topics of far  more versatility and information.

# REFERENCES

[1] Gourav G. Shenoy, Erika H. Dsouza, Sandra Kuebler.Performing Stance Detection on Twitter Data using Computational Linguistics Techniques.

[2] Bartosz Gralewicz.(2018).The TF*IDF Algorithm Explained.

[3] T. Petkovic, Z. Kostanjcar, and P. Pale. E-Mail System for Automatic Hoax Recognition. 2005.

[4] M. Vukovic, K. Pripuzic, and H. Belani. Intelligent Automatic Hoax De- tection System". In: Knowledge-Based and Intelligent Information and En- gineering Systems. Springer, Berlin, Heidelberg, Sept. 2009, pp. 318325.

[5] M. Shari, E. Fink, and J. G. Carbonell. of Internet Scam Using Logistic Regression". In: 2011 IEEE International Conference on Systems, Man, and Cybernetics. Oct. 2011, pp. 21682172.

[6] Brendan Nyhan and Jason Reifler. (2015). The effect of fact-checking on elites: A field experiment on US state legislators. American Journal of Political Science 59(3):628640.

[7] Condren C. (2012). Satire and definition. Humor.

[8] B Adler, L. De Alfaro, S. Mola-Velasco, P. Rosso, and A. West. vandalism detection: Combining natural language, metadata, and reputation features". In: Computational linguistics and intelligent text processing (2011), pp. 277288.

[9] Riffkin, R. (2015, 2015/09/28/). Americans' Trust in Media Remains at Historical Low. Gallup.

[10] Neel Rakholia, Shruti Bhargava.(2015).Is it true? Deep Learning for Stance Detection in News.

[11] I. Augenstein, T. Rocktaschel, A. Vlachos, and K. Bontcheva.(2016). Stance detection with bidirectional conditional encoding. [24]Semeval2016 task 6. [Online].

[12]   M. Del Vicario, A. Bessi, F. Zollo, F. Petroni, A. Scala, G. Caldarelli, H. E.Stanley, and W. Quattrociocchi. Spreading of Misinformation Online". en. In: Proceedings of the National Academy of Sciences 113.3 (Jan. 2016), pp. 554559.

[13]   Hancock J., Woodworth M. Porter, S. (2011). Hungry like a wolf: A word pattern analysis of the language of psychopaths. Legal and Criminological Psychology.

[14]   B Adler, L. De Alfaro, S. Mola-Velasco, P. Rosso, and A. West. vandalism detection: Combining natural language, metadata, and reputation features". In: Computational linguistics and intelligent text processing (2011), pp. 277288.

[15]   Paul Rayson, Andrew Wilson, and Geoffrey Leech. (2001). Grammatical word class variation within the british national corpus sampler. Language and Computers 36(1):295306.

[16]   David B. Buller and Judee K. Burgoon.(1996). Interpersonal deception theory. Communication Theory 6(3):203242.

[17]   Rakhi Joon Archana Singhal.(2017).ANALYSIS OF MWES IN HINDI TEXT USING NLTK.

[18]   M. Potthast, B. Stein, and R. Gerling. vandalism detection in Wikipedia". In: European Conference on Information Retrieval. Springer. 2008, pp. 663668.

[19]   Gregor Leban, Bla Fortuna, Janez Brank Marko Grobelnik.(2014).Event Registry Learning About World Events From News.

[20]   Xiang Zhang, Junbo Zhao, and Yann LeCun. (2015). Character-level convolution networks for text classification. In Advances in Neural Information Processing Systems. pages 649657.

[21]   Marco Bonzanini.(2015).Easy Text Analytics with the Dandelion API and Python.

[22] Rada Mihalcea and Carlo Strapparava.(2009). The lie detector: Explorations in the automatic recognition of deceptive language. In Proceedings of the ACLIJCNLP 2009 Conference Short Papers. Association for Computational Linguistics, pages 309312.

[23] Sepp Hochreiter and Jurgen Schmidhuber.(1997). Long short-term memory. Neural Computation 9(8):17351780.

[24] Pfaff K. L., Gibbs R. W. (1997). Authorial intentions in understanding satirical texts.

[25] Nitin Jindal and Bing Liu. (2008). Opinion spam and analysis. In Proceedings of the 2008 International Conference on Web Search and Data Mining. ACM, pages 219230.

[26] Eunsol Choi, Chenhao Tan, Lillian Lee, Cristian Danescu-NiculescuMizil, and Jennifer Spindel. (2012). Hedge detection as a lens on framing in the gmo debates: A position paper. In Proceedings of the Workshop on Extra-Propositional Aspects of Meaning in Computational Linguistics. Association for Computational Linguistics, pages 7079.

[27] Theresa Wilson, Janyce Wiebe, and Paul Hoffmann. (2005). Recognizing contextual polarity in phrase level sentiment analysis. In Proceedings of Human Language Technology Conference and Conference on Empirical Methods in Natural Language Processing. Association for Computational Linguistics, pages 347354.

[28] Marta Recasens, Cristian Danescu-Niculescu-Mizil, and Dan Jurafsky.(2013). Linguistic models for analyzing and detecting biased language. In Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). Association for Computational Linguistics, pages 16501659.

[29]  Victoria L. Rubin, Yimin Chen, and Niall J. Conroy. (2015). Deception detection for news: Three types of fakes. Proceedings of the Association for Information Science and Technology 52(1):14.

[30]  Giovanni Luca Ciampaglia, Prashant Shiralkar, Luis M. Rocha, Johan Bollen, Filippo Menczer, and Alessandro Flammini. (2015). Computational fact checking from knowledge networks. PLOS ONE 10(6):e0128193.

[31]  Prakhar Biyani, Kostas Tsioutsiouliklis, and John Blackmer.(2016). 8 amazing secrets for getting more clicks: Detecting clickbaits in news streams using article informality. In Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence. AAAI Press, pages 94100.

[32]  William Yang Wang. (2017). Liar, liar pants on fire: A new benchmark dataset for fake news detection. In Proceedings of the Association for Computational Linguistics Short Papers. Association for Computational Linguistics.

[33]  Giovanni Luca Ciampaglia, Prashant Shiralkar, Luis M. Rocha, Johan Bollen, Filippo Menczer, and Alessandro Flammini. (2015). Computational fact checking from knowledge networks. PLOS ONE 10(6):e0128193.

[34]  Matthew L Newman, James W Pennebaker, Diane S Berry, and Jane M Richards. (2003). Lying words: Predicting deception from linguistic styles. Personality and social psychology bulletin 29(5):665675.