

Towards A Headline-based Deception Detection Approach for Data Veracity in Online Digital News

Normala Che Eembi @ Jamil, Iskandar Ishak, Fatimah Sidi, Lilly Suriani Affendey
*Faculty of Computer Science and Information Technology, Universiti Putra Malaysia,
43400 UPM Serdang, Selangor Darul Ehsan, Malaysia
normala.jamil@yahoo.com*

Abstract—Since its existence in 1990's, online news has been the major source of news content for news readers. Unfortunately, based on a number of findings, readers tend to judge on certain event based on the news headlines rather than its contents. With the advancement of mobile and web technologies, it is easier to spread news to others and this unhealthy habits can cause negative impacts towards individuals, organizations or nations that are victimized by the news. This paper proposes a framework to detect deceptive online news based on the news headlines. By having an accurate detection upon deceptive news based on its content, it can assist readers to identify misleading news and to help them to find relevant news for their source of information.

Index Terms—Headlines; Deception; Online News; Framework.

I. INTRODUCTION

One of the characteristics in Big Data is veracity. According to IBM, veracity refers to the level of reliability associated with certain types of data. It also refers to the correctness of the data in Big Data environment such as the uncertain continuous data in GPS sensors or weather conditions or sentiment and truthfulness in human (expressed in social media). Veracity is also a way to find the truthfulness, availability, accountability and authenticity while deception refers to the way of identifying whether verbal expressions are truthful or not as well as whether the overall content is truthful or not [1][2]. One of the issues in data veracity is deception. [3] highlighted deception and truthfulness of data as part of data veracity issues in Big Data which also relates to information quality.

In order to solve the deception of data issues, a number of researches have been conducted to perform detection of deception in texts. Machine learning approaches such as Support Vector Machine (SVM) and Fuzzy technique are among the most popular machine learning approaches used [4]–[7]. SVM holds a major advantage when compared to other machine learning approaches in terms of the theory that can be analyzed in detail and out of the box.

However, the impact of data veracity issues is not just about the content, but it is also about how the information is being retrieved or read. Based on a number of studies, [8][9][10] highlighted that readers tend to only read the headline rather than reading the content and some even straightaway sharing the article without reading the content. Therefore, deceptive headlines may skew the understanding of the reader and they may have misconceptions towards the actual meaning of the content. This is because some of the news used titles which are different in meaning compared to its content in order to catch reader's attention. This brings harms to the masses

when readers especially for those who just read the headlines spread the news based on their understandings of the headlines rather than the news content. Thus, the issue of detecting the deceptive news based on headlines is very important to identify and detect news that contains deceptive element either through its contents or headlines.

II. DECEPTION IN ONLINE NEWS

In 2012, 644 million people have accessed online newspaper sites and this was 42.6% of the total Internet population [11]. This shows that online news is an important source of information for people globally. Any diversion or differences in this news compared to the original news will definitely give impact differently to the people who read it.

News are among the easiest information that people will read and discuss daily but researchers have found out that trustworthiness and the veracity of the news content can be lacking [12][13][14][15]. Similarly, online news shares the same problem but the scale of the problem is greater than the traditional paper based news as online news can be shared easily and spread rapidly through computers and mobile devices regardless of its trustworthiness and veracity.

Based on the literature, the issues of veracity in news portal revolve around three issues; hiding information, adversarial attacks, and masking. The issue of veracity in news reporting always refers to the practice of deceptive writing that sways truthfulness of the news content. Previous studies have reported about hiding information in writing style [16]. In this paper, the researcher used a forensic technique called 'stylometry'. This technique helps to detect authorship of unknown documents. [17] define veracity not only by their written script to identify writers but instead uses content analysis. From their studies, they developed a framework to test the authorship recognition through adversarial attacks. [18] found out that an author could deliberately mask his writing style to give other meanings to the readers that its original content. There are two ways that authors can mask; imitation attack and obfuscation attack. Obfuscation attack is where the authors write the way that it can hide his personal style of writing. Imitation attack is when the author has the ability to write articles that show other authors' style of writing.

Some literature also highlighted the issue of the impact of deceptive writing. [19] reported that when information is hidden with added cognitive for the purpose of deception, changes in human behaviour can be formed. [20] informs in his studies about deceptive information through the media to influence receivers. The researcher wants to know the impact of the media and focuses on computer games. The purpose of

his studies is truthfulness can determine the main sources. This issue will lead to the problem of deception from real news portal.

A number of approaches have been produced in dealing the deceptive writings issue. [21] uses big data technique to protect privacy policies users and used user profiling to collect details of information. Researcher [22] used a framework model based on four elements that are trustor, trustee, trust object and trust part. The trust framework developed for the consumer in the e-commerce domain. Other than that, Van Dam & van de Velden [23] focuses on using a framework to explore and segment user profiles. They profile the Facebook users. Previous studies [24] also develop a web recommendation framework of user profiling based on Probabilistic Latent Semantic Analysis (PLSA) model. The benefit of PLSA is to measure the incidence activities.

A. Deception in News Reporting

News are a packaged information about current events happening in some location [25]. Currently, news is available not just in the traditional paper-based newspaper, but also in the web portal. Web portal-based news has been preferred by news reader rather than the conventional newspaper due to its easiness to access through computers and mobile devices. However, due to the popularity of online news, there are also independent portals which can be developed by either group of peoples or individuals creating online news portals. Other than that, based on [26], content change does exist in news reporting and therefore the veracity of the news must be determined. Deception has been proven to be implemented in news reporting [27][28]. This has caused a number of adverse outcomes such as political unrest, slander and negative perception towards a particular organization, personnel, and country [29]. Hence, it is important for readers to be able to filter out news or portals which are detected as deceptive.

Another problem with deception in news reporting is on the sharing practices. March 2005 [30] shows that the problem where news that has been rebroadcast by news stations did not acknowledge the original source of the news. This may be exploited by some parties or agencies or individuals which may have an interest on inciting rage towards public or war to some extent.

B. Machine Learning-based Deception Detection Approaches

News portal is semi-structured data and data size is large in which a proper approach using the artificial intelligent method as well as data analytic approach must be used. Therefore, machine learning is one of the approaches that has been applied in deception detection domain. Support Vector Machine (SVM) and Bayes technique are among the most popular machine learning approaches used for deception detection. [4]–[7]. For the purpose of comparisons, we choose the most recent researches (starting from 2012) that implemented machine learning technique in text deception detection. Approach by [5] was focused on Chinese Text-based CMC and using Support Vector Machine and produced acquired 90% of accuracy on deception detection.

Another approach by [28] which also implemented Support Vector Machine in deception detection acquired 62.17% of accuracy. Their main focus was on short text on the Twitter dataset. While [31] used the Bayesian technique in deception detection and focus on online review

communities.

For the machine learning technique using fuzzy logic, they focus on employee data to detect deceptive datum with 57.4% of accuracy [32]. Meanwhile [33] have 60-70% of accuracy's result that is focused on cross-cultural content in short essay. Lastly [34] focus on satirical cues content where the domain is satirical and legitimate news with result 90% of accuracy.

Based on previous researches, we found out that all researchers focused their analysis more on the content of the news. However, in terms of reading practices by the masses, studies have shown that many readers only read the news headlines before they judge and act further (such as sharing it with others and judging the actual news) [9][8][10]. [9] also highlighted that headlines are the first impression of news articles and it can drive the way readers perceive the rest of the content associated with them through affecting their way of remembering it. Furthermore, headlines can be disreputably misleading, inaccurate, or ambiguous [35].

Thus, readers struggle to update their memory in order to correct their initial misconceptions. News agencies also regard news headlines as the most crucial part of the information that the readers ought to know [36]. This shows that headline is more important than the rest of the content of the news.

Deceptive headlines can create such a big impact towards the society especially those which are misleading. Among the impact of misleading headlines is the readers tend to be biased towards or away from a specific interpretation misleading headlines can also lead to misconceptions and misinformed behavioural intentions by the readers [9].

Thus, the problem of misleading or deceptive headline can be deemed as data veracity problem. In order to address this problem, the approach for detecting deceptive news based on deceptive headline needs to be promptly addressed. Therefore, we want to propose a deception detection framework that will focus on categorizing deceptive news with misleading or deceptive headline and content to detect deception in online news.

III. PROPOSED FRAMEWORK OF DECEPTION DETECTION

The goal of deception detection is to detect truthfulness of online news based on text features, will be able to determine whether news headline is deceptive. Figure 1 describes the general framework of deception detection for text-based information based on the literature as discussed in this paper.

Our target is to enhance this framework to cater news headline-based deception detection online news as shown in Figure 2.

The framework of the proposed method consists of six main phases. There is dataset, transformation, pre-processing, feature extraction of news headline, implementation of suitable machine learning technique and lastly our proposed deception detection method.

Dataset phase contains the data source which in our case will be news text. The dataset that will be used is the collection of 360 news articles that representatives of the scope US and Canadian national newspaper [34]. It is labelled into two type news which is legitimate and satirical news. Each of news contains 180 of total news. Table 2 shows the detail of dataset.

Table 1
Recent Trend in Deception Detection Using Machine Learning Technique

Authors	Machine learning approaches	Domain	Accuracy	Focus
[34]	SVM	Satirical and Legitimate News	90%	Focus on satirical cues content
[33]	SVM	Short Essay	60-70%	Focus on cross-cultural content
[28]	SVM	Short text	62.17%	Focused short text in social media
[32]	Fuzzy Logic	Employee data	57.4%	Focus on employee data to detect deceptive datum
[5]	SVM	Chinese text	86%	Focused only on Chinese text content
[31]	Bayes	Online review communities data	89.7%	Focus on deceptive online review communities

Table 2
Dataset Details

Details	Legitimate News	Satirical News
Total News	180	180
Total words headline	1744	1937
Average words in each news	9.67	10.76
Average letters in each news	61.62	69.66

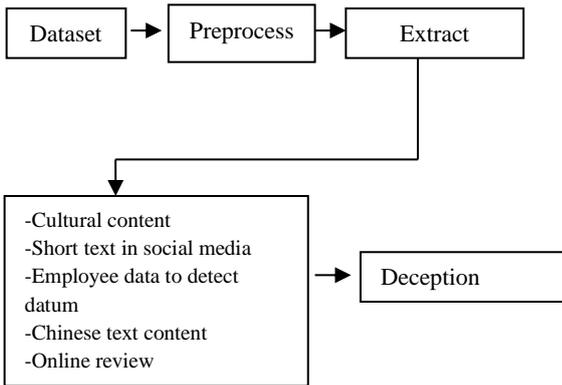


Figure 1: Framework of deception detection

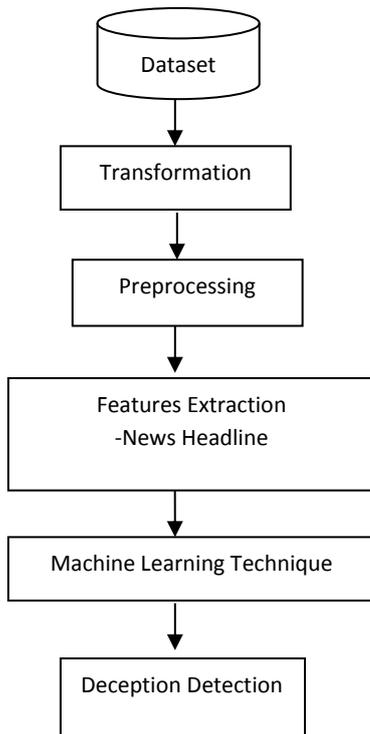


Figure 2: Proposed framework of headline-based deception detection for online news

Transformation phase changes the data into a suitable form that enables processing. The data must be cleaned first as the data comes from any sources. In this case, we will remove all noise and useless information from the data.

After the transformation phase, preprocessing will be conducted to normalize data by removing the unrelated data. The tasks involved in this phase are stop word removal and word stemming. Stop removal word is a technique that is used for removing all the stop words, which is a group of word that has no importance in a sentence. Examples of stop words are ‘the’, ‘and’, ‘to’, ‘a’ and more, while stemming is a technique to find the root of words. For example, words such as ‘waited’, ‘waits’ and ‘waiting’ will be reduced to the root word ‘wait’.

The important and unique part of this framework is the features extraction phase, in which it will be focusing on the news headlines. Then, we will use suitable and relevant machine learning approach to analyze our result. From the analysis, we can predict the deception of an online news headline.

We consider each word in the news headline is very important. Therefore, we are going to use term frequency-inverse document frequency (TF-IDF) method as a part of features extraction in our study. TF-IDF can evaluate how important is a word in the document. Mathematically, TF-IDF is expressed as:

$$TF(t) = \frac{\text{Number of times term } t \text{ appears in document}}{\text{Total number of terms in the document}}$$

$$IDF(t) = \log_e \frac{\text{Total number of documents}}{\text{Number of document with term } t \text{ in it}}$$

After the feature selection process, news headline classification will be determined by the deception detection. The classifier that we will use is Support Vector Machine (SVM). This classification method is the most common method used in research and for news text classification.

IV. EXPERIMENTAL SETUP

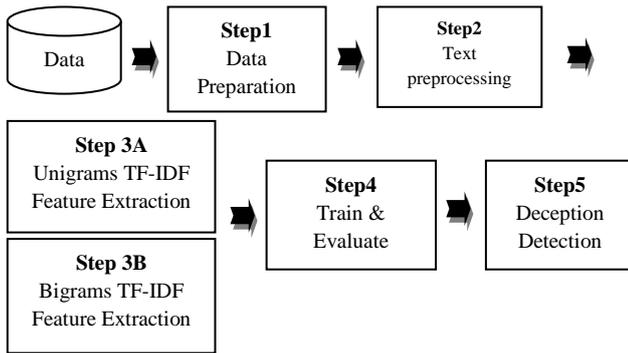


Figure 3: Experimental setup

Figure 3 describes the experimental set up for validating our proposed framework. In step 3A & 3B we define tf-idf using unigrams and bigrams feature extraction. The next step is to train and evaluate the accuracy of the headlines.

Table 3 describes the differences of unigrams and bigrams. Example of the news headline that is used in Figure 4 is “Jeb Bush’s Cerebral Debate Style Faces a Test: Donald Trump” [34].

Table 3
Unigrams and Bigrams for News Headline

Unigrams	Bigrams
('Jeb',)	('Jeb', 'Bush's')
('Bush's',)	('Bush's', 'Cerebral')
('Cerebral',)	('Cerebral', 'Debate')
('Debate',)	('Debate', 'Style')
('Style',)	('Style', 'Faces')
('Faces',)	('Faces', 'a')
('a',)	('a', 'Test:')
('Test:',)	('Test:', 'Donald')
('Donald',)	('Donald', 'Trump')
('Trump',)	

V. FUTURE WORKS AND CONCLUSIONS

Based on the literature, we found out that machine learning techniques such as SVM, Bayesian and Fuzzy are among the popular approaches used by researchers in deception detection on the text. Literature also highlighted that research in deception detection are more focused on the content as a whole rather than the headline. The existence of deceptive headlines in news reporting and the importance of headlines compared to its content were also highlighted by a number of researchers. The impact of deceptive headlines to the masses was also highlighted and we conclude that focusing headlines in terms of deception detection is very important in deceptive deception for online news. Therefore, in this paper, we propose a deception detection framework for online news based on a news headline. As a future work, we will validate and implement our proposed framework to be used as an approach for deception detection for online digital news content.

ACKNOWLEDGMENT

This research was supported by Universiti Putra Malaysia through Putra Grant Scheme-Putra Graduate Initiative (GP-IPS/2016/9478900).

REFERENCES

- [1] Y. R. Tausczik and J. W. Pennebaker, “The Psychological Meaning of Words: LIWC and Computerized Text Analysis Methods,” *J. Lang. Soc. Psychol.*, vol. 29, no. 1, pp. 24–54, 2010.
- [2] V. L. Rubin and T. Vashchilko, “Extending information quality assessment methodology: A new veracity/deception dimension and its measures,” *Proc. Am. Soc. Inf. Sci. Technol.*, vol. 49, no. 1, pp. 1–6, 2012.
- [3] T. Lukoianova and V. L. Rubin, “Veracity roadmap: Is big data objective, truthful and credible?,” *Adv. Classif. Res. Online*, vol. 24, pp. 4–15, 2013.
- [4] M. G. Moghaddam and A. Mustapha, “A Temporal-Focused Trustworthiness to Enhance Trust-based Recommender Systems 20J 3 J 3th International Conference on Intelligent Systems Design and Applications (ISDA),” pp. 219–223, 2013.
- [5] H. Zhang, Z. Fan, J. Zheng, and Q. Liu, “An improving deception detection method in Computer-Mediated Communication,” *J. Networks*, vol. 7, no. 11, pp. 1811–1816, 2012.
- [6] L. Berti-Equille, “Data veracity estimation with ensembling truth discovery methods,” pp. 2628–2636, 2015.
- [7] T. Ekin, F. Leva, F. Ruggeri, and R. Soyer, “Application of bayesian methods in detection of healthcare fraud,” *Chem. Eng. Trans.*, vol. 33, pp. 151–156, 2013.
- [8] D. Dor, “On newspaper headlines as relevance optimizers,” *J. Pragmat.*, vol. 35, no. 5, pp. 695–721, 2003.
- [9] R. Ecker, U.K. Lewandowsky, S. Chang, E.P., Pillai, “The Effects of Subtle Misinformation in News Headlines,” *Uma ética para quantos?*, vol. XXXIII, no. 2, pp. 81–87, 2014.
- [10] D. Q. Wang, “Madness in the Media : Understanding How People With Lived Experience Interpret Newspaper Headlines,” no. April, 2016.
- [11] “Most Read Online Newspapers in the World: Mail Online, New York Times and The Guardian - comScore, Inc.” [Online]. Available: <https://www.comscore.com/Insights/Data-Mine/Most-Read-Online-Newspapers-in-the-World-Mail-Online-New-York-Times-and-The-Guardian>. [Accessed: 24-Jan-2017].
- [12] C. E. Osgood, “Where Do Sentences Come From?,” *Semant. Interdiscip. Read. Philos. Linguist. Psychol.*, pp. 88–105, 1971.
- [13] M. Knapp, R. Hart, and H. Dennis, “An exploration of deception as a communication construct,” *Hum. Commun.*, vol. Fall, no. 1, pp. 15–29, 1974.
- [14] T. Dirsehan and M. Çelik, “Profiling online consumers according to their experiences with a special focus on social dimension,” *Procedia - Soc. Behav. Sci.*, vol. 24, pp. 401–412, 2011.
- [15] M. Kerby and A. Marland, “Media Management in a Small Polity : Political Elites ’ Synchronized Calls to Regional Talk Radio and Attempted Manipulation of Public Opinion Polls,” no. August, 2015.
- [16] S. Afroz, M. Brennan, and R. Greenstadt, “Detecting Hoaxes , Frauds , and Deception in Writing Style Online,” pp. 461–475, 2012.
- [17] J. Wayman, N. Orlans, Q. Hu, F. Goodman, A. Ulrich, and V. Valencia, “Technology Assessment for the State of the Art Biometrics Excellence Roadmap. Volume 2 (of 3). Face, Iris, Ear, Voice, and Handwriter Recognition,” vol. 2, no. JUNE 1987, 2008.
- [18] M. Brennan and R. Greenstadt, “Practical Attacks Against Authorship Recognition Techniques,” *Artif. Intell.*, pp. 60–65, 2009.
- [19] M. G. Frank, M. A. Menasco, and M. O’Sullivan, “Human behavior and deception detection,” *Wiley Handb. Sci. Technol. Homel. Secur.*, 2008.
- [20] H. M. Jung, “Information Manipulation Through the Media,” *J. Media Econ.*, vol. 22, no. 4, pp. 188–210, 2009.
- [21] O. Hasan, B. Habegger, L. Brunie, N. Bennani, and E. Damiani, “A discussion of privacy challenges in user profiling with big data techniques: The EEXCESS use case,” *Proc. - 2013 IEEE Int. Congr. Big Data, BigData 2013*, no. 1, pp. 25–30, 2013.
- [22] W. W. Guo and M. Looi, “A framework of trust-energy balanced procedure for cluster head selection in ireless sensor networks,” *J. Networks*, vol. 7, no. 10, pp. 1592–1599, 2012.
- [23] J.-W. van Dam and M. van de Velden, “Online profiling and clustering of Facebook users,” *Decis. Support Syst.*, vol. 70, pp. 60–72, 2015.
- [24] G. . Xu, Y. . Zhang, and X. . Zhou, “Towards user profiling for web recommendation,” *Lect. Notes Comput. Sci. (including Subser. Lect. Notes Artif. Intell. Lect. Notes Bioinformatics)*, vol. 3809 LNAI, pp. 415–424, 2005.
- [25] R. Nielsen and R. Sambrook, “What is Happening to Television

- News?," *Digit. news Proj. Reuters Inst.*, 2016.
- [26] S. Meraji and C. Tropper, "A Machine Learning Approach for Optimizing," vol. 3, no. 2, 2010.
- [27] K. P. K. Kumar and G. Geethakumari, "Detecting misinformation in online social networks using cognitive psychology," pp. 1–22, 2014.
- [28] V. Pérez-Rosas and R. Mihalcea, "Experiments in Open Domain Deception Detection," 2013.
- [29] Normala, C. Eembi, @ Jamil, I. Bin Ishak, F. Sidi, L. S. Affendey, A. Mamat, N. B. C. E. @ Jamil, I. Bin Ishak, F. Sidi, L. S. Affendey, and A. Mamat, "A Systematic Review on the Profiling of Digital News Portal for Big Data Veracity," *Procedia Comput. Sci.*, vol. 72, pp. 390–397, 2015.
- [30] David Miller, "The age of the fake," *Spin Watch*, 2005.
- [31] M. Ott, C. Cardie, and J. Hancock, "Estimating the prevalence of deception in online review communities," *Proc. 21st Int. Conf. World Wide Web - WWW '12*, pp. 201–210, 2012.
- [32] S. Rajkumar, "Assortment of Uncertainty and Randomness with Fuzzy logic in deception detection for employee database management system using hotchpotch techniques," *World Appl. Sci. J.*, vol. 21, no. 6, pp. 854–857, 2013.
- [33] R. Mihalcea, C. Science, and C. Science, "Cross-cultural Deception Detection," pp. 440–445, 2014.
- [34] V. Rubin, N. J. Conroy, V. L. Rubin, N. J. Conroy, Y. Chen, and S. Cornwell, "Fake News or Truth? Using Satirical Cues to Detect Potentially Misleading News Fake News or Truth? Using Satirical Cues to Detect Potentially Misleading News .," no. April, 2016.
- [35] N. M. Turner, D. G. York, and H. A. Petousis-Harris, "The use and misuse of media headlines: Lessons from the MeNZB??? immunisation campaign," *N. Z. Med. J.*, vol. 122, no. 1291, pp. 22–27, 2009.
- [36] R. A. Metila, "A Discourse Analysis of News Headlines: Diverse Framings for a Hostage-Taking Event," vol. 2, no. 2, pp. 71–78, 2013.