



Volume 22, 2019

ENSEMBLE LEARNING APPROACH FOR CLICKBAIT DETECTION USING ARTICLE HEADLINE FEATURES

Dilip Singh Sisodia

National Institute of Technology
Raipur, Raipur, India

dssisodia.cs@nitrr.ac.in

ABSTRACT

Aim/Purpose	The aim of this paper is to propose an ensemble learners based classification model for classification clickbaits from genuine article headlines.
Background	Clickbaits are online articles with deliberately designed misleading titles for luring more and more readers to open the intended web page. Clickbaits are used to tempted visitors to click on a particular link either to monetize the landing page or to spread the false news for sensationalization. The presence of clickbaits on any news aggregator portal may lead to an unpleasant experience for readers. Therefore, it is essential to distinguish clickbaits from authentic headlines to mitigate their impact on readers' perception.
Methodology	A total of one hundred thousand article headlines are collected from news aggregator sites consists of clickbaits and authentic news headlines. The collected data samples are divided into five training sets of balanced and unbalanced data. The natural language processing techniques are used to extract 19 manually selected features from article headlines.
Contribution	Three ensemble learning techniques including bagging, boosting, and random forests are used to design a classifier model for classifying a given headline into the clickbait or non-clickbait. The performances of learners are evaluated using accuracy, precision, recall, and F-measures.
Findings	It is observed that the random forest classifier detects clickbaits better than the other classifiers with an accuracy of 91.16 %, a total precision, recall, and f-measure of 91 %.
Keywords	authentic news; clickbaits; ensemble learning; natural language processing; performance measures

Accepting Editor Eli Cohen | Received: January 25, 2019 | Revised: March 25, 2019 |
Accepted: March 27, 2019.

Cite as: Sisodia, D. S. (2019). Ensemble learning approach for clickbait detection using article headline features. *Informing Science: the International Journal of an Emerging Transdiscipline*, 22, 31-44. <https://doi.org/10.28945/4279>

(CC BY-NC 4.0) This article is licensed to you under a [Creative Commons Attribution-NonCommercial 4.0 International License](https://creativecommons.org/licenses/by-nc/4.0/). When you copy and redistribute this paper in full or in part, you need to provide proper attribution to it to ensure that others can later locate this work (and to ensure that others do not accuse you of plagiarism). You may (and we encourage you to) adapt, remix, transform, and build upon the material for any non-commercial purposes. This license does not permit you to use this material for commercial purposes.

INTRODUCTION

The availability of news on the internet can be considered both a boon and a bane. As a boon, it provides us with the opportunity to access worldwide happenings just at the click of a button or the press of a key. But, as a bane, the availability of news on the internet has created a frenzy among various websites to earn more income. Clickbaits (Hamblin, 2014; Hoffmann, 2017) are “articles on the internet with content whose primary purpose is to attract attention and encourage visitors to click on a link to a particular webpage” (Zimdars, 2016).

“A clickbait is usually a headline designed to make readers want to click on hyperlinks especially when the links lead to the content of dubious value or interest” (Clickbait, 2018). Clickbaits are mainly used to withhold information intentionally from the readers. They make individuals considerably more inquisitive by giving them something they know a tad bit about, yet not all that much. This leads to a curiosity gap in the minds of the readers. George Loewenstein (1994) has explained this as the information gap theory of curiosity. He defines curiosity as the intrinsic human behavior that is activated when people feel there is a gap between what they know and what they need to know. It is stated in Loewenstein (1994) that theory views curiosity as emerging at the point when consideration gets to be centered around a gap in one’s information. Such information gaps produce the feeling of deprivation called curiosity. The inquisitive individual is inspired to get the missing data to diminish or eliminate the curiosity. Clickbaits (Chen, Conroy, & Rubin, 2015) create this gap in the mind of the readers by using words that tend to lure the reader into opening the link. Some organizations purportedly depend for the most part on clickbait for their traffic.

Most clickbaits fail on their promise of delivering an intriguing story on the web page it points to (Agrawal, 2016), thus, disappointing the reader. Words commanding authority are generally used in clickbaits to assure the reader that the article has to be shown as genuine and not made-up (Blom & Hansen, 2015). Since the headlines consist of words used in such a manner to create an enormous enticing impact on the readers, natural language processing techniques can help in extracting the features of those clickbaits and then we can use classifiers to categorize any given phrase.

The thing with clickbaits is that most people are frequently mindful of this manipulation, but it appears in various forms and requires a lot of efforts. It is stated in Potthast, Köpsel, Stein, and Hagen, (2016) that this has a considerable measure to do with emotion and the part it plays in our daily decision-making processes. Emotional arousal, or the level of physical reaction you have to a feeling, is a key ingredient in clicking behaviors. Another reason stated in Potthast et al. (2016) as to why clickbaits work is the anticipation of pleasure in humans. Reading a clickbait headline, which promises some reward of cute pictures, or anything of that sort, itself creates a sort of pleasure in humans which makes them go on and click on the link. The headline itself is said to give pleasure even before the web page opens.

A clickbait ordinarily has a few of the accompanying qualities as stated in Hurst (2016), an eye-catching and convincing feature, effectively funny or memorable images or videos, humorous tone or offers firmly to a particular emotion, and intended to encourage social sharing.

The remaining of this paper is organized under the following sections. In the next section the related work done for detection of clickbait is discussed. The third section presents the methodology adopted to carry out the present work and described the preparation of training and testing datasets using the different features extracted from article headlines along with the classifiers used in this work. In the Learning Approaches section various measures used for evaluating the performance of learners are described in brief, followed by discussion of the experimental results. Lastly, the paper is concluded with some suggested future work.

RELATED WORK

In this section previously reported work on clickbait detection is discussed in brief. In El-Arini and Tang (2014) an initial survey was conducted to know about the readers' preference regarding the type of contents to read, and it found that 80% of the time individuals favored article headline features that helped them choose if they wanted to peruse the full article before navigating away. In Vijgen (2014) articles with "listicles" are studied. The homogenous structures of the titles of these listicles are similar to the clickbait headlines. It was observed that all the titles contain a cardinal number, and around 85% of all titles begin with these cardinal numbers most popularly ranges between 5 and 25 and an odd number is preferred to an even number. These observations are very helpful in detecting clickbaits. In Reis, Olmo, Prates, Kwak, and An (2015), 69,907 news headlines are used from majors news publishing websites for sentiment analysis of the headlines and discovered that the sentiment of the headlines is unequivocally identified with the prevalence of the news. In Blom & Hansen (2015) porosity in headlines is studied as a way to stimulate interest by analyzing 2000 random headlines collected from a Danish news website. In this study, two forms of forward-references are distinguished, such as discourse deixis (references at discourse level) and cataphora (discourses at phrase level) expressed by demonstrative pronouns, personal pronouns, adverbs, and definite articles, but have not proposed an approach for their detection. A survey on potential methods for automated clickbait detection using textual and non-textual cues is presented in Chen et al. (2015). In Potthast et al. (2016) machine learning models are used to detect clickbaits using extracted features of the tweet headline, the website links and the metadata of the tweet. Biyani, Tsioutsoulouklis, and Blackmer (2016) analyzed properties of clickbait and non-clickbait articles by extracting the headline features, the informality of the web pages, and their URLs. They suggest that informality and forward-reference features lead to the best performance while the performance of all the features combined is better than the performances of individual features. A clickbait detection system based on article and headline pair stance is proposed in Bourgonje, Moreno Schneider, and Rehm (2017). User behavior is incorporated to enhance the performance of clickbait detection. In user behavior analysis model initial clickbait score is calculated using learner and tuned further to improve performance (Zheng, Yao, Jiang, Xia, & Xiao, 2017).

A new webis clickbait corpus is curated using 38,587 annotated tweets and used to evaluate the clickbait detection methods in clickbait challenge 2017 (Potthast et al., 2018). The clickbait issue is also investigated on YouTube videos using metadata of 206K videos. A deep generative variational auto-encoder model was used for classification clickbaits (Zannettou, Chatzis, Papadamou, & Sirivianos, 2018). A new deep learning and metric learning based hybrid techniques integrated with a case based reasoning methodology are proposed for adaptable clickbait detection (López-Sánchez, Herrero, Arrieta, & Corchado, 2018). A deep generative model (Liu, Le, Shu, Wang, & Lee, 2018) is proposed to address the issue of non-availability of large scale labeled data required to train the supervised learning models. In this model, artificial headlines with style transfer are generated from article contents to enlarge the training datasets for performance improvement of clickbait detection. Deep learning models such as recurrent neural networks (RNN) (Anand, Chakraborty, & Park, 2017), Long short term memory (LSTM) (Kumar, Khattar, Gairola, Lal, & Varma, 2018) and convolution neural networks (CNN) (Anand et al., 2017; Zheng et al., 2018) are used to avoid the heavy feature engineering involved in clickbait detection.

This paper proposes an ensemble learners based classification model. Three ensemble learning techniques including bagging, boosting, and random forests are used to design a classifier model for classifying a given headline into the clickbait or non-clickbait.

METHODOLOGY

The methodology adopted in the present work is shown in Figure 1 through a process flow diagram. First, the raw data samples of clickbaits and authentic headlines are collected from various sources.

Then, different features discussed in the literature (Biyani et al., 2016) and potentially useful for detection of clickbaits from datasets of article headlines are decided manually. These features are extracted using programming code and predefined natural language processing libraries and article headlines are represented as feature vectors.

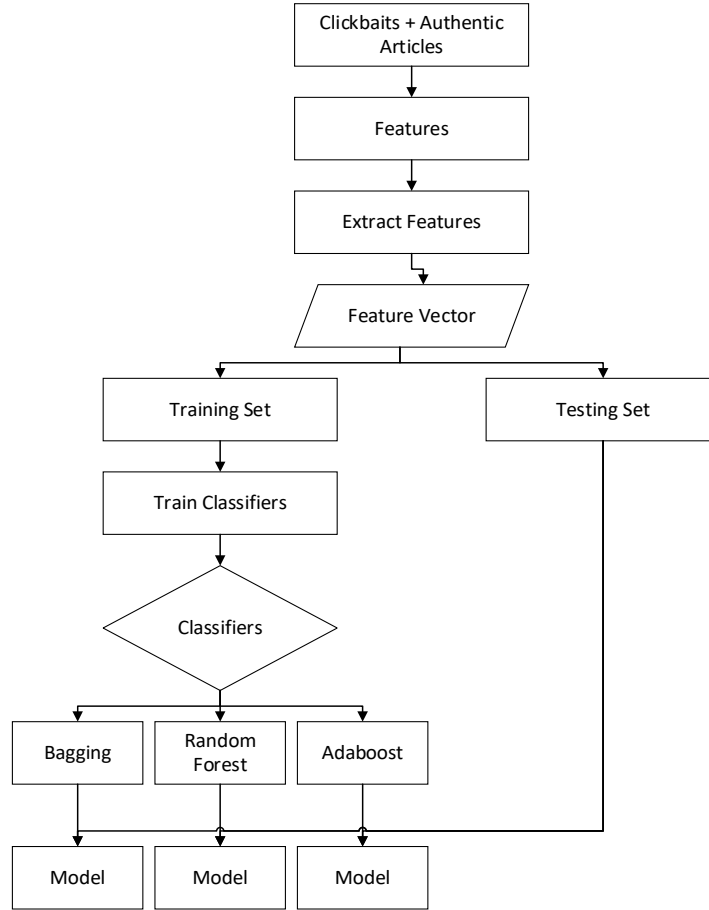


Figure 1. The process flow diagram of an adopted methodology

All the feature vectors of the training samples are passed as input to the classifier along with the class label they belong to, and thus the classifiers undergo learning. After the learning of the classifiers is completed, the feature vectors of the testing samples are sent as input to the models to classify into clickbait or real news headlines.

DATASET DESCRIPTION

The clickbaits used in this paper are the titles of articles published on Buzzfeed – a popular social news and entertainment website (Smith, 2014) – and are obtained from Clickbait Dataset (Retrieved Dec 2014). The entire data consists of around 60,000 clickbaits published in 2014. The authentic news headlines were compiled from Reuters, Associated Press, and The New York Times. They were retrieved using the Article Search API of the New York Times API available for Python (Python, 2017). These headlines were published from January 2013 to March 2016. We have collected around 40,000 news headlines. Table 1 shows the number of samples in each category.

Table 1. Description of Collected No. of Data Samples

Type Data	No. Samples
CLICKBAITS	60,000
GENUINE HEADLINES	40,000

FEATURE SELECTION AND EXTRACTION

Feature selection is a process of selecting most relevant features of data that contribute to the improvement of models to make them less demanding to translate by analysts, shorter training times, and enhanced speculation by lessening overfitting. In the case of clickbaits, the features can be their sentence structures, the presence of certain symbols, etc., which generally distinguishes them from authentic articles. Table 2 contains manually decided 19 attributes (features), which are selected to train the classifiers. All these features are extracted from the article headlines without considering the content of the article as it was used by Biyani et al. (2016).

Table 2. List of Features (Biyani et al., 2016)

S.No	Feature	Description	Type	Example
1	Acronyms	abbreviations	Numeric	<i>25 Awesome DIY Ideas For Bookshelves</i>
2	Adv_adj	adverbs and adjectives	Numeric	<i>Celebrities Riding Invisible Bikes Is Weirdly Hilarious</i>
3	Baity_words	click,happens,next. .etc.	Numeric	<i>Someone Dressed This Dog Up As A Giant Mutant Spider...And What Happens Next Is Hilarious</i>
4	Caps_words	capitalized words	Numeric	<i>CAN'T BE UNSEEN: 30 Art History Snapchats That Are So Inappropriate, But SO Funny</i>
5	Cont_nums	containing numbers (Gardiner, 2015)	Numeric	<i>Sweet Recycling Game Made By An 11-Year-Old</i>
6	Cont_the	containing the	Numeric	<i>The Best Ever Song About Cats</i>
7	Demonstratives	this,that,these,that	Numeric	<i>These 13 Insanely Clever Paint Jobs Will Have You Doing Mind Flips</i>
8	Exclamations	containing ‘!’	Numeric	<i>Check Out All This Cool Vintage Redskins Gear!!!</i>
9	Neg_words	words with negative sentiment	Numeric	<i>New Things Can Be Scary, Even For Corgis</i>
10	Pos_words	words with positive sentiment	Numeric	<i>People Who Think Classic Linkin Park Is Awesome</i>
11	Quoted_words	words within quotes	Numeric	<i>Pepsi Made "Trendy" Clothes In The 80s</i>
12	Question_marks	containing ‘?’	Numeric	<i>Do Games Shape Our Dreams?</i>

S.No	Feature	Description	Type	Example
13	Start_adv	starts with an adverb	Boolean	<i>Beautiful Golden Retriever Puppies</i>
14	Starts_num	starts with a number	Boolean	<i>9 Coffee Swirl Art Masterpieces</i>
15	Swear_words	containing bad words (Finnis, 2015)	Numeric	<i>Disney Put Out An EDM Version Of Let It Go And Its Pretty Damn Good</i>
16	Third_pronouns	he,she,his,...etc.	Numeric	<i>This Woman Has Never Won The Lottery, But She Was Still Able To Ruin Her Own Life</i>
17	Words_5w1h	who,what,when,where,why,how (Mineo, 2017)	Numeric	<i>How Many Hours of Sleep Should You Get?</i>
18	Words_cont_repeated_chars	using letters repeatedly	Numeric	<i>This is sooo funny. Lol!!!</i>
19	Words_title	number of words	Numeric	-

TRAINING SET

The collected data samples are randomly divided into different data sets. Table 3 shows the composition of all the training data sets. The first 3 data sets are balanced, i.e., they contain equal samples of clickbaits and news headlines while the last 2 data sets contain unbalanced data. Five data sets as shown in Table 3 as training inputs to the classifiers, one at a time. These partitions of data are used while testing the machine learning algorithms.

Table 3. The composition of Training Data

Name	Clickbaits	Headlines	Total
DS1	5,000	5,000	10,000
DS2	10,000	10,000	20,000
DS3	15,000	15,000	30,000
DS4	5,000	10,000	15,000
DS5	10,000	5,000	15,000

TESTING SET

We randomly select 10,000 samples each of clickbaits and news headlines to use as test data for the classifiers. This test set is evaluated for all the training partitions of the dataset discussed earlier in Table 3.

ENSEMBLE LEARNING APPROACHES

In this study, three popular ensemble learning techniques such as boosting, bagging, and random forests are evaluated for detection of clickbaits. Ensemble classifier refers to a group of individual classifiers that are cooperatively trained on data set in a supervised classification problem (Rahman & Tasnim, 2014), and ensembles are regularly much more exact than the individual classifiers that make

them up (Dietterich, 2000). It is expressed in Rokach (2010) that the idea of ensemble methodology is to construct a predictive model by incorporating various models. It is understood that ensemble methods can be utilized for enhancing prediction performance. The fundamental thought behind the ensemble methodology is to measure a few individual classifiers and join them to acquire a classifier that outperforms each one of them. In fact, human being tends to look for a few conclusions before making any imperative decision by measuring the individual opinions and consolidate them to achieve our final choice (Polikar, 2006). The brief descriptions of used ensemble classifiers are given as follows.

BAGGING CLASSIFIER

The Bagging is a “bootstrap” ensemble method used to train an individual classifier on an irregular redistribution of the training sets. Each classifier in the ensemble is created with a different random sampling of the training set. Breiman (1996) demonstrated that bagging is viable on “unstable” learning algorithms where little changes in the training set to result in extensive changes in predictions. He claims that the decision tree algorithm is an example of such an unstable learning algorithm. In bagging, the training subsets are drawn arbitrarily (with replacement) from the training set. Homogeneous base classifiers are trained on those subsets. The class picked by most base classifiers is considered to be the final decision of the bagging classifier; every class is picked with equivalent likelihood.

RANDOM FOREST CLASSIFIER

A random forest classifier utilizes a substantial number of individual, unpruned decision trees (Breiman, 1999). Breiman (1999) proposed a random forest, which adds an extra layer of arbitrariness to bagging. This fairly unreasonable technique ends up performing extremely well contrasted with numerous different classifiers and is powerful against overfitting. Likewise, it is exceptionally easy to understand as in it has just two parameters (the quantity of variables in the irregular subset at every node and the number of trees in the forest), and is not extremely touchy to their qualities (Liaw & Wiener, 2002).

ADABOOST CLASSIFIER

AdaBoost (Adaptive Boosting) is a prevalent group algorithm that enhances the basic boosting algorithm using an iterative procedure. Boosting creates data subsets for base classifier training by re-sampling the training patterns, however, by providing the most informative training pattern for each consecutive classifier (Schapire, 1990). In boosting, illustrations that are mistakenly anticipated by past classifiers in the arrangement are picked more frequently than cases that were effectively anticipated (Freund & Schapire, 1996). Friedman, Hastie, Tibshirani (2000) have additionally recommended an optional component that fits together with the expectations of the classifiers as an added substance model utilizing a maximum likelihood rule. As indicated by Quinlan (2006), the fundamental explanation behind AdaBoost's disappointment is overfitting. One conceivable approach to abstain from overfitting is to keep the number of iterations as little as could be expected under the circumstances.

PERFORMANCE EVALUATION MEASURES

When we utilize a classifier model for evaluation, we quite often need to take a gander at the exactness of that model as the number of right forecasts from all expectations made. This is the classifier accuracy. When we have to choose whether it is a sufficient model to take care of the issue, accuracy is by all account, not the only metric for assessing the viability of a classifier. Two other valuable measurements are precision and recall. These two measurements can give much more prominent knowledge into the execution attributes of a classifier.

A false positive (F_p) is a result that indicates a given condition has been fulfilled when it actually has not been fulfilled. A false negative (F_n) is a result which indicates that a condition failed while it was actually successful. True positives (T_p) are relevant items that are correctly identified. True negatives (T_n) are irrelevant items that are correctly identified as irrelevant.

A confusion matrix \mathbf{C} is such that $\mathbf{C}_{i,j}$ is equal to the number of observations known to be in group \mathbf{i} but predicted to be in group \mathbf{j} . A confusion matrix is used to describe the performance of the classifier as shown in Table 4.

Table 4. Confusion Matrix

		Predicted Class	
		Yes	No
Actual Class	Yes	T_p	F_n
	No	F_p	T_n

Accuracy is how close a measured value is to the actual (true) value. It is the proportion of instances whose class the classifier can correctly predict. It can be calculated as shown in Equation (1).

$$Accuracy = \frac{T_p + T_n}{Total\ number\ of\ samples} \quad (1)$$

Precision (P) is defined as the number of true positives over the number of true positives plus the number of false positives. Precision measures the exactness of a classifier. A higher precision implies less false positives, while a lower precision implies more false positives. The value of precision can be calculated using Equation (2).

$$P = \frac{T_p}{T_p + F_p} \quad (2)$$

The recall is defined as the number of true positives over the number of true positives plus the number of false negatives. Recall measures the completeness, or sensitivity, of a classifier. Higher recall implies less false negatives, while lower recall implies more false negatives. The value of recall can be calculated using Equation (3).

$$R = \frac{T_p}{T_p + F_n} \quad (3)$$

Precision and recall can be combined to produce a single metric known as F-measure (F1), which is the weighted harmonic mean of precision and recall. The value of ‘F1’ can be calculated using Equation (4).

$$F1 = \frac{2(P * R)}{P + R} \quad (4)$$

The performances of each classifier are evaluated using these measures.

EXPERIMENTAL RESULTS AND DISCUSSION

Extensive experiments are performed on the dataset described in Table 1. Each of the ensemble classifiers is trained using the 5 data sets individually, and then each classifier is tested using the testing set. In all ensemble learners, the C4.5 decision tree is used as a base learner and Python implementation using Scikit-Learn package (Pedregosa, 2011) are used for experimentation. All the

experiments are performed on a personal computer having 3.40GHz Core i7-4770 with 4.0 GB memory and running under the Microsoft Windows 8.1 Pro.

The default parameters for the ensemble learners are used. For the bagging classifier, the base estimator is the decision tree classifier, and the number of base estimators used is 10. The random forest classifier uses ten estimators, i.e., the number of trees in the forest. For the AdaBoost classifier also, the base estimator used is the decision tree classifier, and the number of estimators, i.e., the maximum number of estimators at which boosting is terminated is 50. To avoid the possibility of overfitting of results experiments are performed using a 10-fold cross-validation scheme.

Figures 2, 3, 4, and 5 respectively show the accuracy, precision, recall, and f-measure of ensemble learners on five datasets. The x-axis represents the different datasets while Y-axis represents the performance metric values (such as accuracy, precision, recall, and f-measure) in percentage.

It is evident from Figure 2 that the accuracy performance of bagging and random forest better and closely related on all five datasets while the accuracy of AdaBoost is consistently very poor on all datasets. Though, the best accuracy performance is recorded for dataset DS3 using random forest ensemble learner model.

It is very difficult to conclude the performance results on the basis of accuracy only. Therefore, other robust performance metrics such as precision, recall, and f-measure are also considered and results reported in Figure 3, 4, and 5 respectively. However, the results of Figure 3, 4, and 5 represent the same performance trend as shown by accuracy measures.

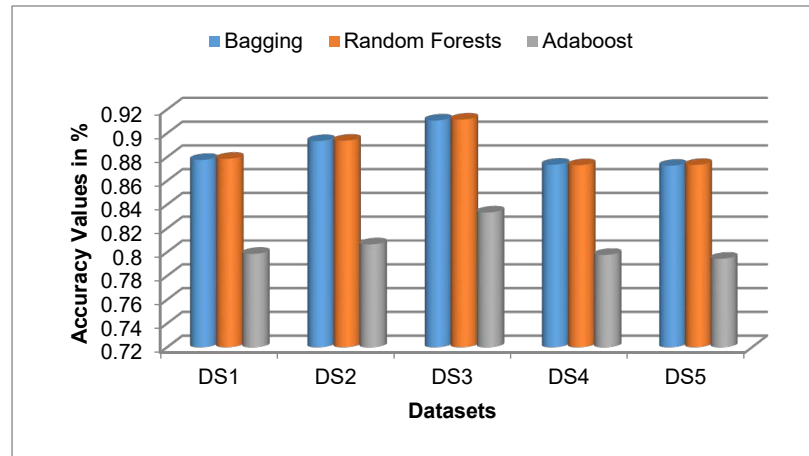


Figure 1. Accuracy of Classifiers

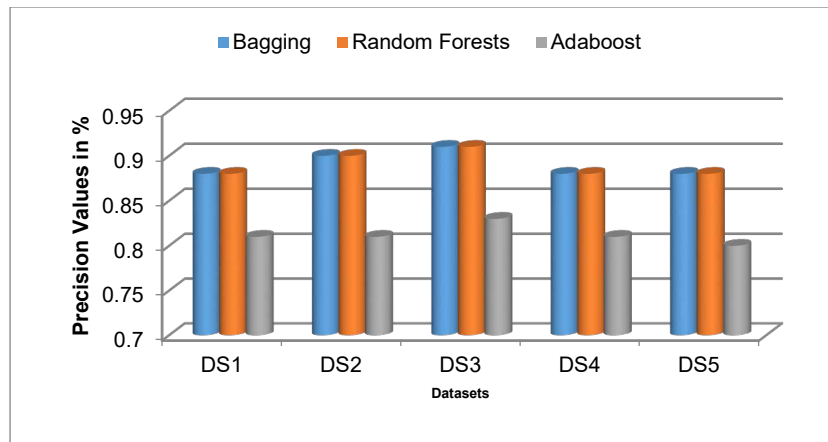


Figure 2. Precision of Classifiers

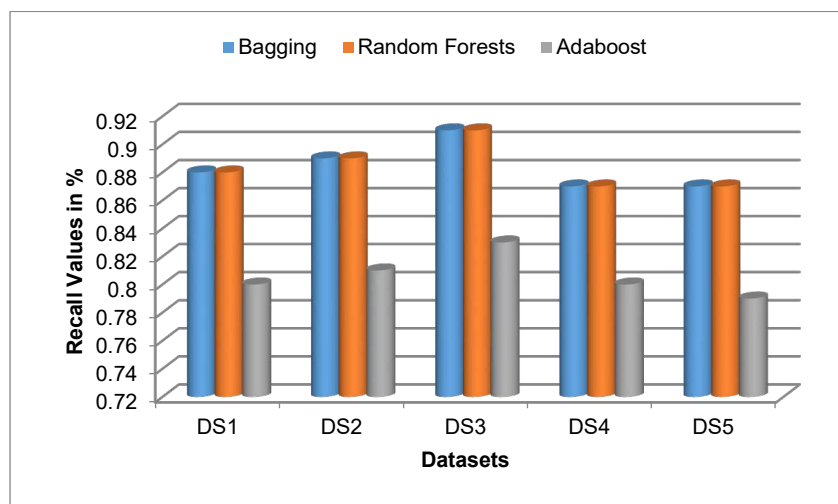


Figure 3. Recall of Classifiers

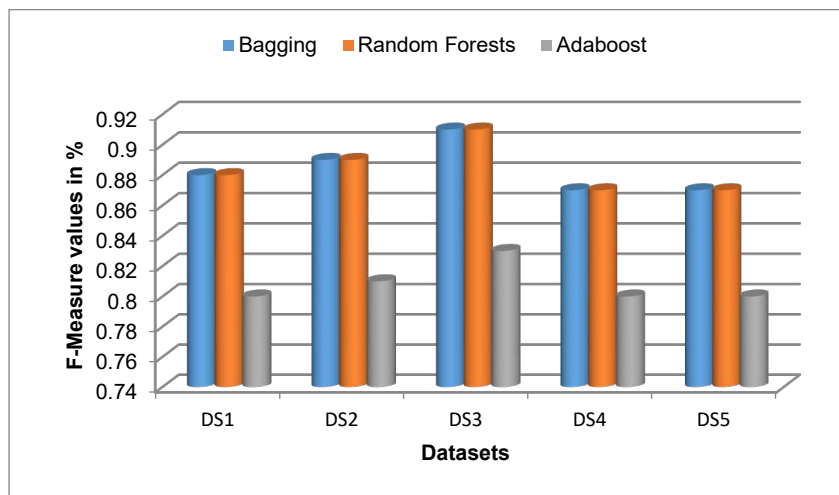


Figure 4. F-Measure values of Classifiers

We can conclude from these figures (Figures 2, 3, 4, and 5) that all the classifiers achieved the highest accuracy, precision, recall and F-measure values for DS3 partition of the dataset. Again, random forest classifier outperforms the rest of the classifiers for all the data partitions, which is closely followed by the bagging classifier. Random forest is giving best performance because it selects a subset of features randomly at each node split of decision tree while bagging consider all features for the same task.

The best performance reported on DS3 data set because DS3 consists of 30,000 samples, which is highest among the used data sets and also contain a balanced data set having an equal number of clickbait (positive) and genuine headlines (negative) samples. All classifiers performed worst on DS4 and DS5 dataset because high imbalance exists in clickbait and genuine headlines samples. This also can be seen from the graphs. DS1 partition also performs relatively lower due to the lesser number of instances.

CONCLUSION

In this paper, the classification clickbaits and authentic news headlines are discussed using ensemble learning algorithms. The dataset of hundred thousand article headlines are collected from news aggregator sites consists of clickbaits and authentic news headlines. The entire data is divided into five different data sets, three of which are balanced and two are unbalanced. The three balanced data sets with a different number of samples of clickbaits and news headlines are used to analyze the performance of the learners. The reason for dividing the data into unbalanced sets is because in the real-world scenario the data available is obviously varying in proportions; sometimes there may be a higher prevalence of clickbaits than news headlines. The features which were considered are those which can be easily extracted from the title of the headline without much difficulty. These features are very simple and fundamental in the process of identifying a clickbait. Accuracy, precision, recall, and f-measure were used as metrics for comparing the performance of the classifiers. It is observed that the balanced data set DS3 with 15000 authentic headlines and 15000 clickbaits gives the best performance for all the classifiers used. The random forest classifier achieves the highest performance with an accuracy of 91.16%. All the metrics show the same trend in the results.

In the present work only features extracted from the title of the headline is used. In the future, many more features can be extracted from other attributes of the headline such as the content (we can consider either the entire content or a part of it), the URL of the link or the website on which the clickbait headline is posted. The content of the clickbait can be cross-checked with the headline to see if it is a clickbait or not. The strategies of natural language processing can be applied in a more advanced manner to detect clickbaits. Sometimes, clickbaits are accompanied by pictures also. The presence of images and the type of image can also be considered as additional features.

REFERENCES

- Agrawal, A. (2016). Clickbait detection using deep learning. In *Next Generation Computing Technologies (NGCT), 2016 2nd International Conference* (pp. 268–272). <https://doi.org/10.1109/ngct.2016.7877426>
- Anand, A., Chakraborty, T., & Park, N. (2017). We used neural networks to detect clickbaits: You won't believe what happened next! In *European Conference on Information* (pp. 541–547). https://doi.org/10.1007/978-3-319-56608-5_46
- Biyani, P., Tsioutsoulklis, K., & Blackmer, J. (2016). “8 amazing secrets for getting more clicks”: Detecting clickbaits in news streams using article informality. *Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence* (pp. 94–100). Retrieved from <https://www.aaai.org/ocs/index.php/AAAI/AAAI16/paper/view/11807/11569>
- Blom, J. N., & Hansen, K. R. (2015). Clickbait: Forward-reference as lure in online news headlines. *Journal of Pragmatics*, 76, 87–100. <https://doi.org/10.1016/j.pragma.2014.11.010>

- Bourgonje, P., Moreno Schneider, J., & Rehm, G. (2017). From clickbait to fake news detection: An approach based on detecting the stance of headlines to articles. In *Proceedings of the 2017 EMNLP Workshop: Natural Language Processing meets Journalism* (pp. 84–89). <http://doi.org/10.18653/v1/w17-4215>
- Breiman, L. (1996). Bagging predictors. *Machine Learning*, 24(2), 123–140.
- Breiman, L. (1999). Random Forests. *Machine Learning*, 45(5), 1–35. <http://doi.org/10.1023/A:1010933404324>
- Chen, Y., Conroy, N. J., & Rubin, V. L. (2015). Misleading online content: Recognizing clickbait as false news. In *Proceedings of the 2015 ACM on Workshop on Multimodal Deception Detection* (pp. 15–19). <https://doi.org/10.1145/2823465.2823467>
- Clickbait. (2018). Definition of Clickbait. In *Merriam-Webster online dictionary*. Retrieved from <https://www.merriam-webster.com/dictionary/clickbait>
- Clickbaits Dataset. (BuzzFeed Articles Retrieved Dec 2014). Retrieved from <https://docs.google.com/spreadsheets/d/1WSx45rT4jZfysmZfzJtjaPO7AxW4XMaJYaCUd5HB2ns/edit#gid=2121187934>
- Dietterich, T. G. (2000). Ensemble methods in machine learning. *Multiple Classifier Systems*, 1857, 1–15. <http://doi.org/10.1007/3-540-45014-9>
- El-Arini, K., & Tang, J. (2014). *Click-baiting*. Facebook Newsroom Retrieved April 19, 2016 from <https://newsroom.fb.com/news/2014/08/news-feed-fyi-click-baiting/>
- Finnis, A. (2015). *List of bad words*. Retrieved from <https://www.buzzfeed.com/alexfinnis/the-100-most-brilliantly-british-swear-words-in-existence>
- Freund, Y., & Schapire, R. R. E. (1996). Experiments with a New Boosting algorithm. *International Conference on Machine Learning*, 148–156. <http://doi.org/10.1.1.133.1040>
- Friedman, J., Hastie, T., & Tibshirani, R. (2000). Additive logistic regression: A statistical view of boosting. *Annals of Statistics*, 28(2), 337–407. <http://doi.org/10.1214/aos/1016218223>
- Gardiner, B. (2015, December 18). You’ll be outraged at how easy it was to get you to click on this headline. *Wired Magazine*, 1–6. Retrieved from <https://www.wired.com/2015/12/psychology-of-clickbait/>
- Hamblin, J. (2014, November 11). It’s everywhere, the clickbait. *The Atlantic*. Retrieved from <https://www.theatlantic.com/entertainment/archive/2014/11/clickbait-what-is/382545/>
- Hoffmann, C. (2017). What is clickbait?(Check all that apply). In C. Hoffmann, *Stupid humanism: Folly as competence in early modern and twenty-first-century culture* (pp. 109–128). Springer. https://doi.org/10.1007/978-3-319-63751-8_5
- Hurst, N. (2016). *To clickbait or not to clickbait? An examination of clickbait headline effects on source credibility*. University of Missouri–Columbia.
- Kumar, V., Khattar, D., Gairola, S., Lal, Y. K., & Varma, V. (2018). Identifying clickbait: A multi-strategy approach using neural networks. In the *41st International ACM SIGIR Conference on Research & Development in Information Retrieval* (pp. 1225–1228). ACM Press. <https://doi.org/10.1145/3209978.3210144>
- Liaw, A., & Wiener, M. (2002). Classification and Regression by randomForest. *R News*, 2(December), 18–22. <http://cogns.northwestern.edu/cbm/LiawAndWiener2002.pdf>
- Liu, H., Le, T., Shu, K., Wang, S., & Lee, D. (2018). Deep headline generation for clickbait detection. In *2018 IEEE International Conference on Data Mining (ICDM)* (pp. 467–476). IEEE. <https://doi.org/10.1109/icdm.2018.00062>
- Loewenstein, G. (1994). The psychology of curiosity: A review and reinterpretation. *Psychological Bulletin*, 116(1), 75–98. <http://doi.org/10.1037/0033-2909.116.1.75>
- López-Sánchez, D., Herrero, J. R., Arrieta, A. G., & Corchado, J. M. (2018). Hybridizing metric learning and case-based reasoning for adaptable clickbait detection. *Applied Intelligence*, 48(9), 2967–2982. <https://doi.org/10.1007/s10489-017-1109-7>

- Mineo, G. (2017). *Why you can't resist clicking on this article: The clickbait conundrum*. Retrieved from <https://blog.hubspot.com/marketing/clickbait-conundrum-history-psychology-ethics>
- Pedregosa. (2011). Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research*, 12, 2825–2830.V
- Polikar, R. (2006). Ensemble-based systems in decision making. *Circuits and Systems Magazine*, 6(3), 21-45. <https://doi.org/10.1109/mcas.2006.1688199>
- Potthast, M., Gollub, T., Komlossy, K., Schuster, S., Wiegmann, M., Fernandez, E. P. G., ... Stein, B. (2018). Crowdsourcing a large corpus of clickbait on Twitter. In *Proceedings of the 27th International Conference on Computational Linguistics* (pp. 1498–1507).
- Potthast, M., Köpsel, S., Stein, B., & Hagen, M. (2016). Clickbait detection. In *European Conference on Information Retrieval* (pp. 810–817). https://doi.org/10.1007/978-3-319-30671-1_72
- Python. (2017). Article Search. Retrieved from <https://pypi.python.org/pypi/NYTimesArticleAPInew/1.0.0>
- Quinlan, J. R. (2006). Bagging, boosting, and C4.5. In *Proceedings of the Thirteenth National Conference on Artificial Intelligence* (Vol. 5, pp. 725–730).
- Rahman, A., & Tasnim, S. (2014). Ensemble classifiers and their applications : A Review. *International Journal of Computer Trends and Technology (IJCTT)*, 10(1), 31–35. <https://doi.org/10.14445/22312803/ijctt-v10p107>
- Reis, J., Olmo, P., Prates, R., Kwak, H., & An, J. (2015). Breaking the news : First impressions matter on online news. *Proceedings of the Ninth International Conference on Weblogs and Social Media*, Oxford, UK, May 26-29, 2015, 357–366.
- Rokach, L. (2010). Ensemble-based classifiers. *Artificial Intelligence Review*, 33(1–2), 1–39. <https://doi.org/10.1007/s10462-009-9124-7>
- Schapire, R. E. (1990). The strength of weak learnability. *Machine Learning*, 5(2), 197–227. <https://doi.org/10.1023/A:1022648800760>
- Smith, B. (2014). Why BuzzFeed doesn't do clickbait. *Buzzfeed*, November, 6.
- Vijgen, B. (2014). The listicle: An exploring research on an interesting shareable new media phenomenon. *Studia Universitatis Babes-Bolyai-Ephemerides*, 59(1), 103-122.
- Zannettou, S., Chatzis, S., Papadamou, K., & Sirivianos, M. (2018). The good, the bad and the bait: Detecting and characterizing clickbait on youtube. In *IEEE Symposium on Security and Privacy Workshops*, SPW 2018 (pp. 63–69). IEEE. <https://doi.org/10.1109/spw.2018.00018>
- Zheng, H. T., Yao, X., Jiang, Y., Xia, S. T., & Xiao, X. (2017, July). Boost clickbait detection based on user behavior analysis. In *Asia-Pacific Web (APWeb) and Web-Age Information Management (WAIM) Joint Conference on Web and Big Data* (pp. 73-80). Springer, Cham. https://doi.org/10.1007/978-3-319-63564-4_6
- Zheng, H. T., Chen, J. Y., Yao, X., Sangaiah, A. K., Jiang, Y., & Zhao, C. Z. (2018). Clickbait convolutional neural network. *Symmetry*, 10(5), 1–12. <https://doi.org/10.3390/sym10050138>
- Zimdars, M. (2016). *False, misleading, clickbait-Y, and satirical "news" sources*. Retrieved from <https://d279m997dpfwgl.cloudfront.net/wp/2016/11/Resource-False-Misleading-Clickbait-y-and-Satirical-%E2%80%9CNews%E2%80%9D-Sources-1.pdf> .

BIOGRAPHY



Dr. Dilip Singh Sisodia received a Ph.D. degree in computer science and engineering from the National Institute of Technology Raipur, India. He earned his Master of Technology and Bachelor of Engineering degrees respectively in information technology (with specialization in artificial intelligence) and computer science & engineering from the Rajiv Gandhi Proudyogiki Vishwavidyalaya (A State Technological University of Madhya Pradesh), Bhopal, India.

Presently, Dr. Sisodia is working as an assistant professor in the department of computer science engineering, National Institute of Technology Raipur. He has over sixteen years of rich academics & research domain experience of various reputed engineering institutes. He has published over 60 refereed articles in journals, conference proceedings and books, published by reputed publishers including IEEE, Springer, Elsevier, SAGE, IOS press, and IGI Global, etc. He is also working as an active reviewer for many international journals and conferences. His current research interests include web usage mining, machine learning, and computational intelligence. Dr. Sisodia is actively associated with various professional societies including IEEE, ACM, CSI, IETE, IE (India), etc.