

Online Learning for Social Spammer Detection on Twitter

Phuc Tri Nguyen · Hideaki Takeda

Received: date / Accepted: date

Abstract Social networking services like Twitter have been playing an import role in people’s daily life since it supports new ways of communicating effectively and sharing information. The advantages of these social network services enable them rapidly growing. However, the rise of social network services is leading to the increase of unwanted, disruptive information from spammers, malware discriminators, and other content polluters. Negative effects of social spammers do not only annoy users, but also lead to financial loss and privacy issues. There are two main challenges of spammer detection on Twitter. Firstly, the data of social network scale with a huge volume of streaming social data. Secondly, spammers continually change their spamming strategy such as changing content patterns or trying to gain social influence, disguise themselves as far as possible. With those challenges, it is hard to directly apply traditional batch learning methods to quickly adapt newly spamming pattern in the high-volume and real-time social media data. We need an anti-spammer system to be able to adjust the learning model when getting a label feedback. Moreover, the data on social media may be unbounded. Then, the system must allow update efficiency model in both computation and memory requirements. Online learning is an ideal solution for this problem. These methods incrementally adapt the learning model with every single feedback and adjust to the changing patterns of spammers overtime. Our experi-

ments demonstrate that an anti-spam system based on online learning approach is efficient in fast changing of spammers comparing with batch learning methods. We also attempt to find the optimal online learning method and study the effectiveness of various feature sets on these online learning methods.

Keywords Social spammer detection · Online learning · Twitter

1 Introduction

Social network services like Twitter have been played a major role in people’s daily life. It offers a convenient and efficient way to communicate and disseminate information. Individuals use Twitter to tweet anything about their concern such as news, jokes, or even their feeling. Companies and organizations use Twitter as an effective channel to connect with their customers, promote or sell their products. With these advantages, Twitter has been increasingly used for large-scale information dissemination in various fields of human life such as marketing, journalism or public relations.

Nevertheless, the popularity of Twitter has led to the rise of unwanted, disruptive information from social spammers. Twitter spammers (Wikipedia, 2016) are defined as malicious users who try to gain social influence and generate spamming contents which negatively impact on legitimate users. Spammers are motivated to launch various of attacks such as stealing personal information of users (Bilge et al, 2009), spreading viruses, malware (Grier et al, 2010), phishing attacks, or compromise suspicious fake followers. Social spammers do not only annoy users, but also lead to financial loss and

Phuc Tri Nguyen
University of Information Technology,
Ho Chi Minh, Vietnam
E-mail: phucnt@uit.edu.vn

Hideaki Takeda
National Institute of Informatics
Tokyo, Japan
E-mail: takeda@nii.ac.jp

privacy issues of users. Therefore, the problem of social spamming is a serious issue prevalent on Twitter. Characterizing and detecting social spammers can keep Twitter as a spam-free environment and improve the quality of user experiences.

There are two main challenges of spammer detection on Twitter.

1. The first challenge is how to process an enormous amount of Twitter data. Today, Twitter services handle more than 2.8 billion requests and store 4.5 petabytes of time series data every minutes (Twitter, 2016). We need an approach that can be able to scale up to handle a huge volume of data with limited computation capacity.
2. Fast change of spamming patterns is the second challenge. The social spammer detection usually seems like a endless game between spammers and anti-spam systems. Spammers continually change their spamming strategy to fool the anti-spam systems. An approach that can be able to adapt to the complex, and fast changing of data is needed.

There were some previous studies on the Twitter spammer detection for years. Their approaches address this problem as the task of classifying a Twitter user into a spammer or a legitimate user. By analytic the spammers behaviors, they proposed effective features which related to content-based and network-based of users and then built a traditional batch learning model to detect spammers for future data. However, such batch learning models are less efficient due to rapid changing and quick evolution of spammers. Moreover, because of limited resources, it is very expensive if we gather all spamming patterns from batches to train the learning model.

One efficient approach for the fast evolve and large-scale of social spammer detection on Twitter is online learning. The online learning method continually updates the existing model while data arrives, as opposed to batch learning techniques which learn from the entire data. In this paper, we study how to apply online learning on the social spammer detection on Twitter. Our experiment on Honey Pot dataset (Lee et al, 2011) and 1KS-10KN (Yang et al, 2012) dataset indicate the effectiveness of the online learning approach for this problem.

The main contributions of our work are outlined as follows:

1. Successfully apply online learning for the problem of social spammer detection on Twitter. Our ex-

periments show that online learning approaches efficiently reflect with the fast changing of data.

2. Find an optimal online learning method for social spammer detection on Twitter. We evaluated 16 online learning algorithms and find the Soft Confidence-Weight algorithm achieves the best performance.
3. Evaluate the effectiveness of four different feature sets when applying online learning on this problem. The best result was observed by the combination of 2 sets of user network features and user activities features. This results indicated that user profile features and user content features are less robust than user network features and user activities features.

The remaining of this report is organized as follows. In the next section, we will present a brief overview of related works on spammer detection. A description of the methodology applied will be described in Section 3. Then an experimental study showing the effectiveness of online learning is presented in Section 4. Finally, we conclude and discuss future work in Section 5.

2 Related Works

2.1 Spammer detection in other platforms

Spammers have been around us since the beginning of the electronic communication and adapted through the development of technology. Spam detection problem is a serious issue, and it has been studied for years on various platforms such as SMS (Gómez Hidalgo et al, 2006), email (Blanzieri and Bryl, 2008), and the Web (Webb, 2006). A popular and well-developed approach for spam detection is based on machine learning techniques. They extracted effective features from historical data and built a supervised learning model. This model will be used to classify new data as either spam or legitimate user/message.

2.2 Spammer detection on Twitter

As a result of the popularity of social media like Twitter, spammers are turning into the fast growing in this platform. There were some previous studies to tackle this problem.

Some studies focus on analyzing the spammer characteristics on Twitter. (Yardi et al, 2009) explored the behavior of spammers from the entire life cycle of #robot-pickupline hashtag. According to their observations, spammers tend to send more messages and network interaction with others. Thus, the higher ratio of followers to followings of a user have, the higher probability that

the user is a spammer. (Grier et al, 2010) studied perspective of spammers and click-through behaviors. Additionally, they had already evaluated the effectiveness of backlists to prevent spamming. (Thomas et al, 2011b) collected the suspended accounts in 7 months from August 2010 to March 2011 and then they studied the characteristic of spammer account, tweet behavior, and spam campaign. (Ghosh et al, 2012) focused on the link farm on Twitter by analyzing the suspended accounts. They observed that most of the link farms came from new users. They also proposed Collusion Rank to demote the ranking of link farms on Twitter.

(Benevenuto et al, 2010) addressed the study of spammer detection on Twitter trending topics. They used an SVM classifier to distinguish between spammers and legitimate users based on the basic of tweet contents and user profile information. Using the social honeypot to collect spammers is an interested work studied by (Lee et al, 2010, 2011). After analysis spammer behavior, they extracted user profile features and user network features and built a supervised learning classifier to identify spammers. However, the approach requires a lot of time for observation spamming evidence. Additionally, the collected data is often biased because it was only received content polluted from active spammers who was following the honeypot accounts.

(Yang et al, 2011) observed that the proposed features from previous works were less effective with the evolving of spammers. They utilized ten new features for spammer detection on Twitter and evaluated these features with the existing ones. Their experiments indicated that using their new features give a better result for spammer detection problem. (Ferrara et al, 2014) used the information related to tweet content, user network, sentiment, and temporal patterns of activity for detecting Twitter bots. Some approaches have assessed the safety or suspiciousness of URLs in tweets as a mean to identify spam tweets (Thomas et al, 2011a; Cao and Caverlee, 2015; Wang et al, 2013; Lee and Kim, 2012).

More recent work has investigated the relationship between automation and spamming. In the study of (Amleshwaram et al, 2013), a system for automated spammers detection is described. Features related to automation have been exploited to adapt to the changing structure of Twitters spammers population. An analysis of automated activity on Twitter was presented on (Chu et al, 2010), and a system that detects the automation of an account is described in (Chu et al, 2012).

Previous works on spam detection on Twitter can be

summarized as follows: collecting and analyzing spammer behaviors, define and extracting effective features, using supervised learning algorithms to build the statistical classifier to detect spammers. However, the behavior of spammers changes too fast in social media. It is hard for the batch-learning system to adapt to the evolving of spammers.

3 Online Learning for Social Spam Detection

To detect spammer on Twitter, we propose a framework based on online supervised classification method. The classify model will be incrementally updated real-time. The overall framework is presented in the Fig 1. In the training step, given m users with their identity label, the system will extract the features vector of each user. The detail of the extracted features will be talking in the feature represent spammers or legitimate users. Feature vectors and identity labels will be used to built the classifier model H_m . Given one more user, the feature extraction module will be used to extract the features of the user and update classifier model H_m to H_{m+1} .

3.1 Features represent spammers or legitimate users

A spammer and a legitimate user have motivation differently in posting tweets or doing social activities. We can assume that the characteristics of spammers are quite different to legitimate users. The features to present a user include profile, social networks, activities and content of posted tweets.

3.1.1 User profile features

The profile features are extracted from a user's Twitter profile and consist of:

- the length of screen name
- whether the user profile has a description, the length of description
- whether the user profile has a URL
- the longevity (age) of an account (hours, days, weeks)

3.1.2 User networks features

The following features are used to characterize a user's social networks which mainly extracted from friendship information.

- number of users following (friends)
- number of followers
- the ratio of number of following to number of followers

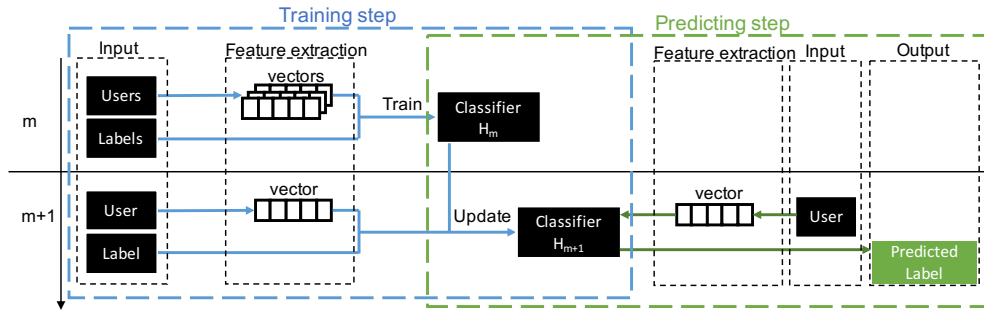


Fig. 1: The online learning framework for spammer detection on Twitter

- the reputation of a user which is calculated by the ratio between the number of followers and the total number of followers and following. User's reputation $= \frac{followers}{following+followers}$
- the following rate which is calculated the ratio between number of following and the longevity of account (hours, days, weeks)
- the followers rate which is calculated the ratio between number of followers and the longevity of account (hours, days, weeks)
- number of bidirectional friend $following \cap follows$
- the percentage of bidirectional friends with following $\frac{following \cap follows}{following}$
- the percentage of bidirectional friends with followers $\frac{following \cap follows}{follows}$
- the standard deviation of followings
- the standard deviation of followers

3.1.3 User activity features

These features category capture user's social activities such as posted tweets or retweets.

- number of posted tweets
- number of posted tweets per hours, days, weeks
- number of content similarity of posted tweets by a user.
- number of direct mentions (e.g., @username) per posted tweet
- number of direct mentions (e.g., @username) per hours, days, weeks
- number of URLs per tweet
- number of URLs per hours, days, weeks
- number of hashtags per tweet
- number of hashtags per hours, days, weeks
- number of retweets per tweet
- number of retweets per hours, days, weeks

3.1.4 User content features

Capture the linguistic properties of the text of each tweet such as part-of-speech tagging, the number of spam words from spamming words dictionary¹, Linguistic Inquiry and Word Count (LIWC) and sentiment features.

- Part-of-speech tagging provides the syntactic information of a sentence and has been used in the natural language processing for measuring text informativeness. In detail implementation, we use the Twitter-specific tagger (Gimpel et al, 2011).
- Number of spam words is generated by matching a famous list of spam words. This list contains over 21,000 phrases, patterns, and keywords commonly used by spammers and comment bots in usernames, email addresses, link text. Since the masking behavior can dramatically decrease the proportion of spam tweets in a spamming account, applying this feature on an account content may not be helpful in detecting complex spamming accounts.
- LIWC dictionary is used to analyze text statistically and find psychologically-meaningful categories (Pennebaker et al, 2001). There are 68 defined categories in LIWC dictionary. We use LIWC dictionary to compute 68 user's personality features. These features may help to determine the personality of spammers or legitimate users.
- In psychological, the micro expressions (Matsumoto and Hwang, 2011) play a distinct role in detecting deception. Inspired by this work, I explore the sentiment information could help capture deceptions of spammers. In this work, I use the list of lexicons from (Dodds et al, 2011) for generating the sentiment features.

¹ <https://github.com/splorp/wordpress-comment-blacklist>

3.2 Online learning algorithms

This section briefly presents the online learning algorithms which we use for our evaluation. Informal, the online learning algorithms are trying to solve an online classification problem over a sequence of pairs $\{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}$, where each x_m is an example's feature vector and $y_m \in \{0, 1\}$ is its label. At each step m during training, the algorithm makes a label prediction $h_m(x_m)$, which for linear classifier is $h_m(x) = \text{sign}(w_m x)$.

After making a prediction, the algorithm receives the actual label y_m . Then the algorithms compute the loss $l(y_m, \hat{y}_m)$ based on some criterion to measure the difference between the prediction and the revealed true label y_m . The learner finally decides when and how to update the classification model at the end of each learning step bases on the result of the loss function.

In this paper, we have no vested interest in any particular strategy for online learning. We simply focus on the application of online learning on the problem social spammer detection on Twitter.

In our experiment, we use the LIBOL an online learning tool which was developed by (Hoi et al, 2014). This tool consists of existing state-of-the-art online learning algorithms for large-scale online classification tasks. In details, these online learning algorithms can be grouped into two following categories: first-order online learning algorithms and second-order online learning algorithms.

First-order online learning algorithms include the algorithms that only keep updating one classification function. The examples algorithms in this categories are following:

- Perceptron: the classical online learning algorithm (Rosenblatt, 1958)
- ALMA: the Approximate Maximal Margin algorithm (Gentile, 2002)
- ROMMA: the Relaxed Online Maximum Margin algorithm (Li and Long, 2002)
- OGD: the Online Gradient Descent algorithm (Zinkevich, 2003)
- PA: the Passive Aggressive algorithm (Crammer et al, 2006)

Second-order online learning algorithms have been explored in recent years. The major family of this categories assume the weight vector follows a Gaussian distribution $w \sim N(\mu, \Sigma)$ with mean vector $\mu \in \mathbb{R}^d$, covariance matrix $\Sigma \in \mathbb{R}^{d \times d}$ and dimensional vector

Table 1: Statistics of 2 two datasets: Honey Pot (Lee et al, 2011) and 1KS-10KN (Yang et al, 2012)

	<i>Spammers</i>		<i>Legitimate Users</i>	
	# Users	# Tweets	# Users	# Tweets
Honey Pot	22,223	2,353,473	19,276	3,259,693
1KS-10KN	1,000	145,095	10,000	1,209,521

space d . The examples algorithms in this categories are following:

- SOP: the Second Order Perceptron algorithm (Cesa-Bianchi et al, 2005)
- CW: the Confidence-Weight learning algorithm (Crammer et al, 2009a)
- IELLIP: the online learning algorithms by improved ellipsoid method (Yang et al, 2009)
- ARROW: the Adaptive Regularization of Weight Vectors algorithm (Crammer et al, 2009b)
- NARROW: the New variant of Adaptive Regularization (Orabona and Crammer, 2010)
- NHERD: the Normal Herding method via Gaussian Herding (Crammer and Lee, 2010)
- SCW: the Soft Confidence Weight algorithms (Wang et al, 2012)

4 Experiments

In this section, we conduct the experiments to evaluate the effectiveness of online learning over the social spammer detection on Twitter. To demonstrate the effectiveness of online learning, we address the following questions:

1. What is the accuracy of online learning methods compare with the batch learning methods on the single dataset?
2. Do online learning algorithms provide a better adaptation on the data distribution changing?
3. Which online algorithms are most appropriate for our application?
4. How about the effectiveness of four different feature sets when applying online learning on the social spammer detection on Twitter?

4.1 Datasets

In our experiments, we use two Twitter datasets. The statistics of two datasets are presented in Table 1.

4.1.1 Social Honeypot Dataset:

(Lee et al, 2011) created 60 Tweeter accounts to attract spammers. After seven months, from December 2009 to August 2010, his team collected the information of 41,499 users and 5,613,166 their tweets. 22,223 users were labeled as spammers and 19,276 legitimate users. The ratio of spammers to legitimate users is around 1:1.

4.1.2 1KS-10KN Dataset:

This dataset was collected by (Yang et al, 2012) from April 2010 to July 2010. This dataset contains the information of 11,000 users and 1,354,616 their tweets. They labeled 1,000 users as spammers and 10,000 legitimate users. The ratio of spammers to legitimate in the 1KS-10KN data is 1:10.

In two dataset, we don't have the information about the changing of network information of users through a specific time interval. We only have a snapshot of network information at the ending of collected data. For our work, we adapted the datasets with our purpose by randomly splitting each dataset to 20 parts and considering each part as a time interval.

4.2 Effectiveness of Online learning on single dataset

In this section, we start by evaluating the effectiveness of online learning over the batch learning method regarding classification cumulative error rate on Honey Pot dataset and 1KS-10KN dataset. A lower the cumulative error rate means the better performance. In details, we compare two online learning algorithms - SCW (Wang et al, 2012) and ALMA (Gentile, 2002) against two different training set configurations of the batch learning - Random Forest algorithms. We choose the batch learning - Random Forest algorithm because it was produced the highest performance on previous researches: (Lee et al, 2011; Yang et al, 2011). 2 online learning algorithms: SCW (Wang et al, 2012) and ALMA (Gentile, 2002) were chosen because they gave the lowest and second-lowest cumulative error rate in our experiments which will be presented in Table 4. Assuming that our system can only process one part of the dataset, we conduct experiments to every single part of data.

Figure 2 and Table 2 shows the classification cumulative error rates for online learning method and batch learning method on Honey Pot dataset and 1KS-10KN dataset. The x-axis shows the percentage of the dataset

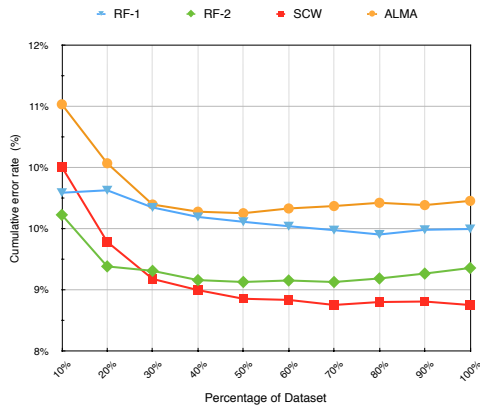
Table 2: Effectiveness of online learning on single dataset

%	HoneyPot				1KS-10KN			
	RF1	RF-2	SCW	ALMA	RF-1	RF-2	SCW	ALMA
10	0.101	0.098	0.104	0.112	0.018	0.020	0.038	0.055
15	0.102	0.095	0.099	0.108	0.014	0.018	0.030	0.053
20	0.101	0.091	0.094	0.105	0.013	0.015	0.025	0.048
25	0.099	0.090	0.092	0.102	0.012	0.014	0.021	0.044
30	0.099	0.091	0.089	0.099	0.013	0.015	0.019	0.044
35	0.099	0.090	0.090	0.099	0.014	0.015	0.019	0.044
40	0.098	0.089	0.088	0.098	0.014	0.016	0.018	0.044
45	0.098	0.090	0.088	0.099	0.013	0.015	0.017	0.042
50	0.097	0.089	0.087	0.098	0.013	0.015	0.016	0.040
55	0.097	0.089	0.087	0.099	0.012	0.014	0.016	0.039
60	0.096	0.089	0.087	0.099	0.012	0.014	0.015	0.037
65	0.097	0.090	0.087	0.099	0.012	0.014	0.014	0.035
70	0.096	0.089	0.086	0.099	0.012	0.014	0.014	0.033
75	0.095	0.089	0.086	0.099	0.012	0.015	0.014	0.032
80	0.095	0.090	0.086	0.099	0.012	0.015	0.014	0.032
85	0.096	0.090	0.086	0.099	0.013	0.015	0.014	0.031
90	0.096	0.090	0.086	0.099	0.013	0.015	0.014	0.030
95	0.096	0.090	0.086	0.099	0.013	0.015	0.014	0.029
100	0.096	0.091	0.086	0.100	0.013	0.015	0.014	0.029

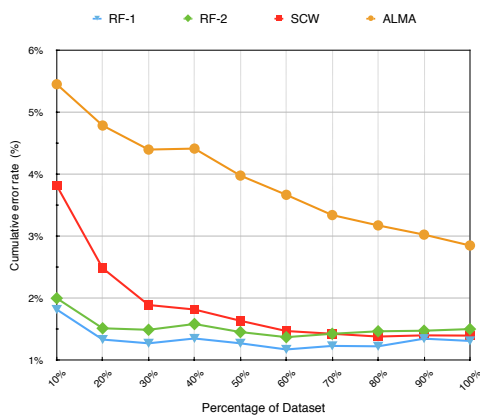
and the y-axis shows the cumulative error rate: percentage of miss-classified examples for all user up to this part. The RF-1 represent using the batch learning - Random Forest to train once on the first part of 20 splitting data. Then the model will be used to test on all the remaining parts. The RF-2 uses the same setting with RF-1 but retrain the model for every interval. For example, RF-2 train on part 2 of dataset and test on part 3. SCW and ALMA are 2 online learning used to make a single update over a cumulative training data.

In our experiment on single dataset, the distribution of data does not change much since we do not have the real streaming data on social spammer detection. We try to adapt the data Honey Pot and 1KS-10KN with our purpose. At the beginning part of data, the cumulative error rate of two online learning algorithms: SCW and ALMA are lower than batch learning: RF-1 and RF-2. However, the next results on the remaining parts show that SCW and ALMA faster reduce the cumulative error rate than batch learning. It means that online learning gives a more rapid adaptability when testing with other parts of data.

These results tested on single dataset indicate that online learning can give a comparable result with batch-learning when the distribution of data does not change much. Although online learning does not achieve the better result compare with batch learning techniques, it allows faster adapt to new data.



(a) Honey Pot dataset



(b) 1KS-10KN dataset

Fig. 2: Effectiveness of online learning on single dataset

4.3 Effectiveness of Online learning on combine dataset

In the real world problem, social spammers always change their spamming strategy; it makes rapid changing in the distribution of data. It explained why we study the ability to adapt to the changing distribution when to combine two datasets. We use the same online learning algorithms and two setting of batch learning in the previous section and evaluate on the combined dataset. In particular, this dataset contains 20 first parts from Honey Pot dataset and 20 next parts from 1KS-10KN dataset.

Figure 3 show the result of social spammer detection on the combine of Honey Pot and 1KS-10KN dataset. In the RF-1, the cumulative error rate becomes significant increase when testing on data from 1KS-10KN. It is indicated the classifier model fail to detect spammers when data distribution changes. If we want to achieve better accuracy, we need to retrain the model with the

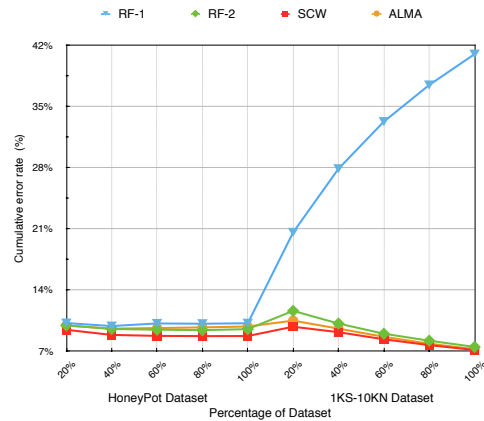


Fig. 3: Effectiveness of Online learning on combine dataset

new data from the 1KS-10KN dataset. In RF-2, the classifier will be re-trained from the prior interval, gave the better performance than RF-1. It means that the new data eventually can help to reduce the cumulative error rate. However, online learning algorithms - SCW and ALMA gave the better adaptability with the changing of data distribution since it gave the lowest cumulative error rate compare with the batch-learning method. The experiment results show that online learning very appropriate with the fast changing of spammers.

4.4 Which online algorithms are most appropriate for social spammer detection

In this section, we evaluate which of the online learning algorithms are the best suited to the social spammer detection. We tested 16 online learning algorithms implemented using the LIBOL tool. All of the experiments run on the combined dataset discussed in the previous section.

Table 4 shown the comparison of the performance of online learning algorithms. We find that the cumulative error rate ranged from 0.265 to 0.071. The Soft Confidence-Weight showed the best performance compare with the other online learning algorithms.

4.5 Effectiveness of various feature set combinations

In this section, we trained 2 online learning (SCW-the highest performance and ALMA-the second-highest performance) with different feature sets (e.g. UP: User Profile features, UN: User Network features, UA: User Activities features, UC: User Content features) and various combination of the feature sets on the mixture of

Table 4: Performance comparison of online learning algorithms

	%	Perceptron	ROMMA	aROMMA	ALMA	OGD	PA	PA1	PA2	SOP	IELLIP	CW	NHERD	AROW	NAROW	SCW	SCW2
Honey Pot	10%	0.137	0.164	0.168	0.107	0.095	0.145	0.126	0.131	0.137	0.165	0.148	0.097	0.093	0.103	0.105	0.104
	15%	0.129	0.150	0.155	0.102	0.096	0.143	0.128	0.132	0.131	0.159	0.146	0.097	0.092	0.104	0.098	0.101
	20%	0.125	0.148	0.154	0.099	0.092	0.140	0.127	0.130	0.127	0.156	0.147	0.093	0.089	0.102	0.094	0.098
	25%	0.126	0.150	0.155	0.098	0.090	0.137	0.126	0.128	0.123	0.154	0.146	0.090	0.088	0.102	0.093	0.099
	30%	0.123	0.151	0.156	0.096	0.088	0.137	0.125	0.126	0.120	0.154	0.145	0.087	0.085	0.098	0.091	0.097
	35%	0.124	0.152	0.158	0.097	0.089	0.138	0.126	0.129	0.121	0.155	0.150	0.088	0.086	0.098	0.091	0.097
	40%	0.122	0.150	0.158	0.096	0.087	0.135	0.123	0.125	0.118	0.152	0.149	0.087	0.085	0.097	0.089	0.095
	45%	0.122	0.151	0.159	0.096	0.087	0.136	0.124	0.126	0.118	0.153	0.151	0.087	0.085	0.097	0.089	0.095
	50%	0.120	0.150	0.158	0.096	0.086	0.136	0.123	0.125	0.118	0.152	0.152	0.086	0.084	0.096	0.087	0.094
	55%	0.121	0.152	0.159	0.096	0.085	0.136	0.124	0.125	0.118	0.153	0.154	0.086	0.084	0.096	0.087	0.094
	60%	0.121	0.153	0.160	0.096	0.086	0.135	0.123	0.125	0.117	0.153	0.154	0.086	0.084	0.096	0.087	0.093
	65%	0.121	0.154	0.162	0.097	0.086	0.136	0.123	0.125	0.118	0.154	0.155	0.086	0.084	0.096	0.087	0.093
	70%	0.121	0.155	0.163	0.097	0.085	0.136	0.123	0.125	0.118	0.154	0.156	0.086	0.083	0.095	0.087	0.093
	75%	0.122	0.155	0.162	0.097	0.085	0.135	0.123	0.125	0.119	0.154	0.156	0.085	0.083	0.095	0.087	0.093
	80%	0.122	0.157	0.162	0.097	0.086	0.135	0.122	0.125	0.120	0.153	0.157	0.086	0.084	0.095	0.087	0.094
	85%	0.122	0.157	0.163	0.097	0.086	0.135	0.123	0.125	0.120	0.154	0.158	0.086	0.084	0.095	0.087	0.094
	90%	0.123	0.158	0.164	0.098	0.086	0.136	0.123	0.125	0.121	0.155	0.159	0.086	0.084	0.095	0.087	0.094
95%	0.123	0.158	0.164	0.098	0.085	0.136	0.123	0.125	0.121	0.155	0.159	0.086	0.084	0.095	0.087	0.093	
100%	0.123	0.159	0.165	0.098	0.086	0.136	0.123	0.126	0.121	0.155	0.161	0.086	0.085	0.095	0.087	0.094	
1KS - 10KN	5%	0.129	0.162	0.166	0.109	0.102	0.135	0.123	0.125	0.124	0.155	0.155	0.100	0.096	0.111	0.096	0.110
	10%	0.134	0.164	0.168	0.108	0.116	0.133	0.120	0.123	0.129	0.153	0.151	0.110	0.104	0.125	0.098	0.125
	15%	0.135	0.165	0.169	0.108	0.128	0.130	0.117	0.120	0.130	0.151	0.147	0.122	0.116	0.138	0.098	0.140
	20%	0.134	0.166	0.170	0.105	0.135	0.128	0.113	0.118	0.132	0.148	0.144	0.130	0.124	0.148	0.098	0.151
	25%	0.134	0.166	0.168	0.103	0.141	0.124	0.110	0.115	0.131	0.144	0.140	0.138	0.130	0.158	0.097	0.163
	30%	0.135	0.166	0.167	0.100	0.147	0.122	0.107	0.112	0.130	0.142	0.137	0.146	0.136	0.168	0.095	0.175
	35%	0.132	0.170	0.170	0.098	0.152	0.120	0.104	0.110	0.131	0.141	0.135	0.153	0.142	0.177	0.094	0.184
	40%	0.130	0.170	0.169	0.096	0.156	0.118	0.101	0.109	0.130	0.139	0.132	0.158	0.147	0.185	0.092	0.193
	45%	0.129	0.171	0.169	0.093	0.158	0.116	0.099	0.107	0.128	0.138	0.131	0.163	0.150	0.193	0.089	0.200
	50%	0.128	0.171	0.169	0.091	0.162	0.114	0.096	0.105	0.129	0.135	0.129	0.170	0.156	0.202	0.088	0.209
	55%	0.129	0.174	0.171	0.089	0.165	0.114	0.095	0.105	0.128	0.136	0.128	0.175	0.160	0.209	0.086	0.217
	60%	0.126	0.172	0.169	0.086	0.166	0.112	0.092	0.103	0.126	0.134	0.127	0.179	0.163	0.216	0.084	0.224
	65%	0.124	0.172	0.168	0.084	0.167	0.110	0.090	0.100	0.124	0.131	0.124	0.182	0.165	0.221	0.082	0.229
	70%	0.122	0.173	0.169	0.082	0.167	0.110	0.088	0.100	0.123	0.132	0.124	0.185	0.166	0.227	0.080	0.235
	75%	0.121	0.174	0.169	0.080	0.168	0.109	0.086	0.099	0.123	0.131	0.123	0.189	0.169	0.233	0.079	0.241
	80%	0.119	0.176	0.170	0.079	0.168	0.108	0.085	0.098	0.122	0.130	0.121	0.192	0.171	0.240	0.077	0.246
	85%	0.119	0.178	0.173	0.077	0.168	0.108	0.084	0.098	0.122	0.130	0.121	0.196	0.173	0.245	0.076	0.251
90%	0.119	0.177	0.174	0.075	0.168	0.107	0.083	0.097	0.121	0.129	0.120	0.200	0.175	0.251	0.074	0.257	
95%	0.117	0.178	0.175	0.074	0.167	0.106	0.081	0.096	0.119	0.129	0.118	0.202	0.176	0.256	0.073	0.261	
100%	0.115	0.180	0.176	0.072	0.166	0.107	0.080	0.097	0.118	0.129	0.118	0.204	0.176	0.262	0.071	0.265	

Honey Pot dataset and 1KS-10KN dataset. The experiment results are reported in Figure 4, Table 5 and Table 6.

The result of social spammer detection by using SCW algorithm and ALMA algorithm are quite similar. The result of user profile features and user content features lower than user network features and user activity features when using a single feature set. It is consistent because social spammers are easy to fake their profile as a normal user on Twitter. Moreover, spammers quickly change their content information. It makes the user content features become noisy and less efficient.

In all of the experiments, the result of the combination of user network features and user activity features gives the best result in term of cumulative error rate 0.068 by using SCW algorithm and 0.056 by using ALMA. These results indicate that this two kind of features set are robustness with the changing spamming patterns and stable across the choice of online learning algorithms.

5 Conclusion and Future Work

The social spammer detection on Twitter is sophisticated and adaptable to game the system by continually change their content and network patterns. To handle fast evolving social spammers, we suggest using the online learning incrementally update classifier model when the newly spamming pattern occurred. Our experiment results show that the approach is effective in dynamically changing spamming strategies of spammers comparing with other batch learning method. Additionally, we address that the online learning - Soft Confident-Weight achieve the best result compare with other online algorithms. We also studied the effectiveness of four feature sets on two online learning algorithms - SCW and ALMA. The experiments show that user network features and user activities features are more robustness than user profile features and user content features and stable across online learning algorithms.

In near future, the amount of available data has risen

Table 5: Evaluation of various feature set combinations on the SCW algorithm

	%	UP	UN	UA	UC	UP+UN	UP+UA	UP+UC	UN+UA	UN+UC	UA+UC	UP+UN+UA	UP+UN+UC	UN+UA+UC	UP+UN+UC+UA
Honey Pot	10%	0.354	0.152	0.214	0.188	0.133	0.204	0.189	0.103	0.110	0.188	0.096	0.109	0.103	0.104
	15%	0.365	0.150	0.204	0.184	0.137	0.195	0.182	0.100	0.109	0.184	0.094	0.110	0.099	0.100
	20%	0.360	0.147	0.203	0.183	0.133	0.193	0.179	0.097	0.105	0.182	0.091	0.104	0.094	0.095
	25%	0.362	0.146	0.204	0.184	0.132	0.193	0.177	0.096	0.102	0.182	0.091	0.102	0.093	0.093
	30%	0.359	0.145	0.202	0.182	0.130	0.191	0.175	0.093	0.101	0.181	0.089	0.100	0.090	0.091
	35%	0.360	0.146	0.201	0.182	0.131	0.192	0.174	0.094	0.101	0.180	0.090	0.099	0.091	0.090
	40%	0.358	0.146	0.199	0.179	0.130	0.191	0.172	0.093	0.098	0.177	0.090	0.097	0.089	0.088
	45%	0.359	0.145	0.201	0.179	0.130	0.192	0.172	0.093	0.098	0.178	0.090	0.097	0.088	0.088
	50%	0.357	0.143	0.200	0.178	0.128	0.191	0.171	0.092	0.097	0.176	0.089	0.096	0.087	0.087
	55%	0.358	0.143	0.199	0.177	0.128	0.191	0.170	0.092	0.096	0.176	0.088	0.095	0.087	0.087
	60%	0.358	0.143	0.199	0.176	0.127	0.190	0.168	0.091	0.096	0.175	0.088	0.095	0.087	0.087
	65%	0.359	0.143	0.199	0.176	0.128	0.190	0.168	0.091	0.096	0.175	0.089	0.095	0.087	0.087
	70%	0.359	0.143	0.198	0.176	0.127	0.190	0.168	0.091	0.096	0.175	0.088	0.094	0.086	0.086
	75%	0.360	0.142	0.199	0.177	0.127	0.190	0.169	0.091	0.096	0.176	0.088	0.094	0.086	0.086
	80%	0.361	0.142	0.199	0.177	0.127	0.191	0.169	0.091	0.096	0.176	0.088	0.094	0.087	0.087
	85%	0.359	0.142	0.198	0.176	0.127	0.189	0.168	0.091	0.096	0.175	0.089	0.095	0.086	0.086
90%	0.359	0.142	0.198	0.176	0.127	0.189	0.167	0.091	0.096	0.175	0.089	0.095	0.086	0.087	
95%	0.360	0.143	0.197	0.176	0.128	0.189	0.167	0.091	0.096	0.175	0.088	0.095	0.086	0.086	
100%	0.360	0.143	0.198	0.175	0.128	0.189	0.168	0.091	0.096	0.175	0.089	0.095	0.086	0.086	
1KS-10KN	5%	0.369	0.138	0.206	0.196	0.132	0.202	0.186	0.098	0.101	0.191	0.097	0.101	0.098	0.099
	10%	0.376	0.133	0.202	0.211	0.127	0.208	0.204	0.097	0.100	0.203	0.105	0.099	0.106	0.107
	15%	0.381	0.128	0.197	0.223	0.122	0.212	0.219	0.096	0.097	0.212	0.112	0.097	0.111	0.113
	20%	0.387	0.123	0.192	0.234	0.118	0.212	0.235	0.096	0.095	0.218	0.117	0.095	0.114	0.116
	25%	0.391	0.119	0.186	0.243	0.113	0.207	0.249	0.094	0.093	0.223	0.122	0.093	0.115	0.117
	30%	0.396	0.114	0.181	0.249	0.109	0.202	0.260	0.093	0.090	0.226	0.126	0.091	0.115	0.118
	35%	0.400	0.111	0.176	0.254	0.106	0.196	0.271	0.092	0.089	0.229	0.130	0.089	0.116	0.119
	40%	0.405	0.107	0.170	0.260	0.102	0.190	0.281	0.090	0.087	0.231	0.133	0.087	0.116	0.119
	45%	0.410	0.104	0.165	0.266	0.099	0.184	0.290	0.087	0.085	0.232	0.136	0.086	0.115	0.118
	50%	0.411	0.101	0.160	0.270	0.096	0.178	0.299	0.085	0.084	0.233	0.139	0.084	0.115	0.119
	55%	0.414	0.098	0.155	0.274	0.094	0.173	0.306	0.083	0.082	0.234	0.142	0.083	0.115	0.118
	60%	0.417	0.095	0.151	0.277	0.091	0.168	0.313	0.081	0.081	0.234	0.145	0.081	0.114	0.117
	65%	0.421	0.092	0.147	0.279	0.089	0.163	0.321	0.079	0.079	0.233	0.147	0.080	0.113	0.116
	70%	0.424	0.090	0.143	0.281	0.087	0.159	0.329	0.078	0.078	0.232	0.150	0.079	0.112	0.115
	75%	0.426	0.088	0.139	0.284	0.085	0.155	0.335	0.076	0.076	0.231	0.152	0.077	0.112	0.114
	80%	0.428	0.086	0.136	0.286	0.082	0.151	0.340	0.074	0.076	0.231	0.154	0.076	0.111	0.114
85%	0.430	0.084	0.132	0.288	0.081	0.147	0.345	0.072	0.075	0.230	0.157	0.075	0.111	0.113	
90%	0.432	0.082	0.129	0.289	0.079	0.143	0.349	0.071	0.073	0.229	0.159	0.074	0.111	0.113	
95%	0.434	0.080	0.126	0.290	0.078	0.140	0.355	0.069	0.072	0.228	0.161	0.073	0.110	0.112	
100%	0.435	0.079	0.123	0.290	0.076	0.137	0.361	0.068	0.071	0.227	0.162	0.072	0.109	0.111	

steadily. It imposes a computational burden on the single system. In consequences, we need an approach can be able to perform in a distributed fashion. With this motivation, the future works will be focused on how to build the scalable distributed spammer detection system.

References

- Amleshwaram AA, Reddy N, Yadav S, Gu G, Yang C (2013) Cats: Characterizing automation of twitter spammers. In: Communication Systems and Networks (COMSNETS), 2013 Fifth International Conference on, IEEE, pp 1–10
- Benevenuto F, Magno G, Rodrigues T, Almeida V (2010) Detecting spammers on twitter. In: In Collaboration, Electronic messaging, Anti-Abuse and Spam Conference (CEAS)
- Bilge L, Strufe T, Balzarotti D, Kirda E (2009) All your contacts are belong to us: Automated identity theft attacks on social networks. In: Proceedings of the 18th International Conference on World Wide Web, ACM, New York, NY, USA, WWW '09, pp 551–560, DOI 10.1145/1526709.1526784, URL <http://doi.acm.org/10.1145/1526709.1526784>
- Blanzieri E, Bryl A (2008) A survey of learning-based techniques of email spam filtering. *Artif Intell Rev* 29(1):63–92, DOI 10.1007/s10462-009-9109-6, URL <http://dx.doi.org/10.1007/s10462-009-9109-6>
- Cao C, Caverlee J (2015) Detecting spam urls in social media via behavioral analysis. In: Advances in Information Retrieval, Springer, pp 703–714
- Cesa-Bianchi N, Conconi A, Gentile C (2005) A second-order perceptron algorithm. *SIAM Journal on Computing* 34(3):640–668
- Chu Z, Gianvecchio S, Wang H, Jajodia S (2010) Who is tweeting on twitter: human, bot, or cyborg? In: Proceedings of the 26th annual computer security applications conference, ACM, pp 21–30

Table 6: Evaluation of various feature set combinations on the ALMA algorithm

	%														
		UP	UN	UA	UC	UP+UN	UP+UA	UP+UC	UN+UA	UN+UC	UA+UC	UP+UN+UA	UP+UN+UC	UN+UA+UC	UP+UN+UC+UA
Honey Pot	10	0.359	0.137	0.227	0.205	0.132	0.207	0.201	0.101	0.126	0.197	0.096	0.122	0.103	0.107
	15	0.362	0.140	0.216	0.203	0.135	0.196	0.194	0.103	0.123	0.193	0.096	0.121	0.101	0.105
	20	0.360	0.135	0.213	0.202	0.131	0.196	0.192	0.100	0.120	0.194	0.095	0.117	0.098	0.102
	25	0.360	0.134	0.213	0.200	0.131	0.196	0.190	0.100	0.118	0.193	0.095	0.114	0.097	0.100
	30	0.359	0.131	0.213	0.198	0.128	0.194	0.186	0.097	0.115	0.192	0.091	0.111	0.094	0.096
	35	0.360	0.132	0.212	0.197	0.129	0.196	0.185	0.099	0.113	0.191	0.093	0.111	0.095	0.096
	40	0.359	0.132	0.210	0.193	0.127	0.194	0.180	0.098	0.111	0.188	0.091	0.108	0.093	0.095
	45	0.360	0.132	0.211	0.196	0.127	0.197	0.183	0.098	0.111	0.189	0.093	0.108	0.094	0.095
	50	0.357	0.130	0.211	0.194	0.126	0.196	0.182	0.097	0.110	0.188	0.092	0.107	0.094	0.095
	55	0.359	0.131	0.211	0.194	0.126	0.197	0.183	0.097	0.109	0.188	0.091	0.107	0.094	0.095
	60	0.359	0.130	0.210	0.193	0.126	0.196	0.181	0.096	0.109	0.187	0.092	0.107	0.094	0.095
	65	0.360	0.130	0.211	0.193	0.126	0.196	0.182	0.096	0.109	0.187	0.092	0.108	0.094	0.095
	70	0.360	0.130	0.210	0.193	0.125	0.196	0.182	0.096	0.108	0.187	0.092	0.108	0.094	0.094
	75	0.359	0.129	0.211	0.193	0.125	0.196	0.182	0.096	0.109	0.188	0.092	0.108	0.094	0.095
	80	0.359	0.130	0.212	0.193	0.125	0.196	0.182	0.097	0.109	0.187	0.092	0.108	0.095	0.095
	85	0.357	0.130	0.211	0.192	0.125	0.195	0.180	0.097	0.108	0.187	0.093	0.108	0.095	0.095
90	0.356	0.130	0.210	0.192	0.125	0.195	0.180	0.097	0.109	0.187	0.092	0.109	0.095	0.095	
95	0.357	0.130	0.210	0.193	0.125	0.195	0.180	0.098	0.109	0.187	0.093	0.109	0.095	0.096	
100	0.357	0.131	0.210	0.193	0.125	0.195	0.181	0.098	0.110	0.187	0.093	0.109	0.095	0.096	
1KS-10KN	5	0.364	0.144	0.203	0.198	0.139	0.198	0.191	0.100	0.120	0.189	0.102	0.120	0.106	0.105
	10	0.370	0.152	0.194	0.198	0.149	0.191	0.193	0.096	0.122	0.186	0.100	0.124	0.107	0.105
	15	0.375	0.158	0.185	0.197	0.156	0.183	0.193	0.092	0.122	0.181	0.097	0.125	0.106	0.103
	20	0.379	0.160	0.178	0.196	0.160	0.175	0.192	0.088	0.122	0.176	0.095	0.126	0.104	0.101
	25	0.381	0.161	0.171	0.194	0.162	0.168	0.190	0.085	0.121	0.171	0.093	0.125	0.102	0.099
	30	0.383	0.159	0.164	0.191	0.161	0.162	0.188	0.082	0.120	0.166	0.091	0.125	0.099	0.097
	35	0.385	0.157	0.158	0.189	0.160	0.157	0.186	0.079	0.120	0.161	0.086	0.124	0.097	0.094
	40	0.385	0.153	0.153	0.186	0.157	0.151	0.184	0.077	0.118	0.156	0.083	0.123	0.095	0.092
	45	0.387	0.148	0.148	0.184	0.152	0.146	0.182	0.074	0.117	0.152	0.080	0.122	0.092	0.090
	50	0.389	0.143	0.143	0.182	0.148	0.141	0.180	0.072	0.116	0.147	0.079	0.121	0.090	0.088
	55	0.391	0.139	0.138	0.181	0.144	0.137	0.179	0.070	0.115	0.143	0.077	0.120	0.088	0.085
	60	0.391	0.135	0.134	0.178	0.140	0.133	0.177	0.068	0.113	0.139	0.075	0.119	0.085	0.083
	65	0.391	0.132	0.130	0.176	0.137	0.129	0.175	0.066	0.113	0.135	0.073	0.118	0.083	0.081
	70	0.390	0.129	0.127	0.174	0.134	0.125	0.173	0.064	0.112	0.132	0.071	0.117	0.081	0.079
	75	0.390	0.126	0.123	0.172	0.131	0.122	0.171	0.063	0.111	0.129	0.067	0.116	0.079	0.077
	80	0.390	0.123	0.120	0.170	0.128	0.119	0.169	0.061	0.110	0.126	0.065	0.115	0.077	0.076
85	0.390	0.120	0.117	0.169	0.125	0.116	0.168	0.060	0.110	0.123	0.063	0.114	0.076	0.074	
90	0.389	0.118	0.114	0.168	0.123	0.113	0.166	0.059	0.109	0.120	0.062	0.113	0.074	0.073	
95	0.389	0.115	0.112	0.166	0.120	0.110	0.165	0.058	0.108	0.117	0.060	0.112	0.073	0.071	
100	0.386	0.113	0.109	0.165	0.117	0.108	0.163	0.056	0.107	0.114	0.059	0.111	0.071	0.070	

Chu Z, Gianvecchio S, Wang H, Jajodia S (2012) Detecting automation of twitter accounts: Are you a human, bot, or cyborg? Dependable and Secure Computing, IEEE Transactions on 9(6):811–824

Crammer K, Lee DD (2010) Learning via gaussian herding. In: Advances in neural information processing systems, pp 451–459

Crammer K, Dekel O, Keshet J, Shalev-Shwartz S, Singer Y (2006) Online passive-aggressive algorithms. The Journal of Machine Learning Research 7:551–585

Crammer K, Dredze M, Pereira F (2009a) Exact convex confidence-weighted learning. In: Advances in Neural Information Processing Systems, pp 345–352

Crammer K, Kulesza A, Dredze M (2009b) Adaptive regularization of weight vectors. In: Advances in neural information processing systems, pp 414–422

Dodds PS, Harris KD, Kloumann IM, Bliss CA, Danforth CM (2011) Temporal patterns of happiness and information in a global social network: Hedonomet-

rics and twitter. PloS one 6(12):e26,752

Ferrara E, Varol O, Davis C, Menczer F, Flammini A (2014) The rise of social bots. arXiv preprint arXiv:14075225

Gentile C (2002) A new approximate maximal margin classification algorithm. The Journal of Machine Learning Research 2:213–242

Ghosh S, Viswanath B, Kooti F, Sharma NK, Korlam G, Benevenuto F, Ganguly N, Gummadi KP (2012) Understanding and combating link farming in the twitter social network. In: Proceedings of the 21st International Conference on World Wide Web, ACM, New York, NY, USA, WWW '12, pp 61–70, DOI 10.1145/2187836.2187846, URL <http://doi.acm.org/10.1145/2187836.2187846>

Gimpel K, Schneider N, O'Connor B, Das D, Mills D, Eisenstein J, Heilman M, Yogatama D, Flanigan J, Smith NA (2011) Part-of-speech tagging for twitter: Annotation, features, and experiments. In: Proceedings of the 49th Annual Meeting of the Asso-

Table 3: Effectiveness of Online learning on combine dataset

	%	RF-1	RF-2	SCW	ALMA
Honey Pot	10	0.102	0.101	0.105	0.107
	15	0.103	0.102	0.098	0.102
	20	0.102	0.100	0.094	0.099
	25	0.102	0.096	0.093	0.098
	30	0.100	0.096	0.091	0.096
	35	0.101	0.097	0.091	0.097
	40	0.099	0.095	0.089	0.096
	45	0.101	0.097	0.089	0.096
	50	0.102	0.096	0.087	0.096
	55	0.102	0.095	0.087	0.096
	60	0.102	0.095	0.087	0.096
	65	0.102	0.095	0.087	0.097
	70	0.102	0.094	0.087	0.097
	75	0.102	0.094	0.087	0.097
	80	0.101	0.094	0.087	0.097
85	0.101	0.095	0.087	0.097	
90	0.101	0.095	0.087	0.098	
95	0.101	0.095	0.087	0.098	
100	0.102	0.095	0.087	0.098	
IKS - 10KN	5	0.130	0.131	0.096	0.109
	10	0.157	0.126	0.098	0.108
	15	0.182	0.120	0.098	0.108
	20	0.206	0.116	0.098	0.105
	25	0.227	0.112	0.097	0.103
	30	0.246	0.108	0.095	0.100
	35	0.263	0.105	0.094	0.098
	40	0.279	0.102	0.092	0.096
	45	0.294	0.099	0.089	0.093
	50	0.308	0.096	0.088	0.091
	55	0.321	0.093	0.086	0.089
	60	0.333	0.090	0.084	0.086
	65	0.344	0.088	0.082	0.084
	70	0.355	0.086	0.080	0.082
	75	0.365	0.084	0.079	0.080
80	0.375	0.082	0.077	0.079	
85	0.384	0.080	0.076	0.077	
90	0.393	0.078	0.074	0.075	
95	0.402	0.077	0.073	0.074	
100	0.410	0.075	0.071	0.072	

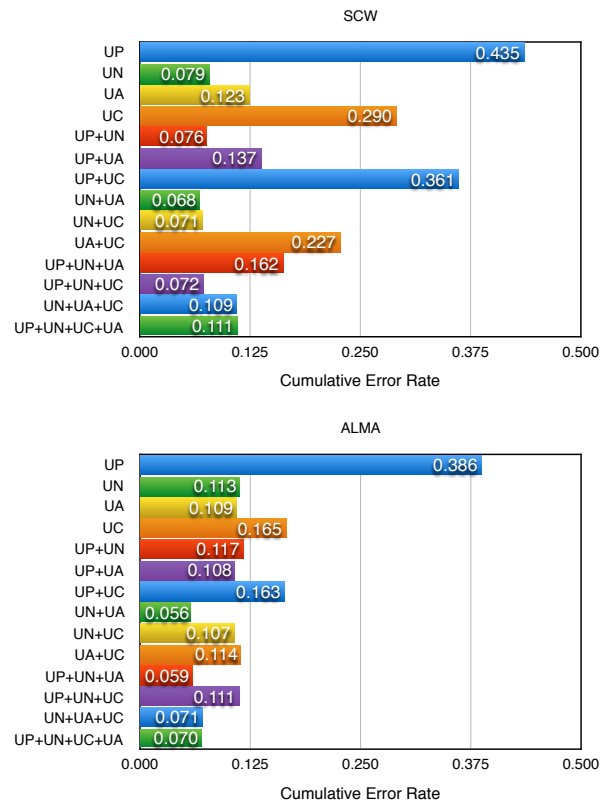


Fig. 4: Evaluation of various feature set combination on SCW algorithm and ALMA algorithm

ciation for Computational Linguistics: Human Language Technologies: Short Papers - Volume 2, Association for Computational Linguistics, Stroudsburg, PA, USA, HLT '11, pp 42–47, URL <http://dl.acm.org/citation.cfm?id=2002736.2002747>

Gómez Hidalgo JM, Bringas GC, Sáenz EP, García FC (2006) Content based sms spam filtering. In: Proceedings of the 2006 ACM Symposium on Document Engineering, ACM, New York, NY, USA, DocEng '06, pp 107–114, DOI 10.1145/1166160.1166191, URL <http://doi.acm.org/10.1145/1166160.1166191>

Grier C, Thomas K, Paxson V, Zhang M (2010) @ spam: the underground on 140 characters or less. In: Proceedings of the 17th ACM conference on Computer and communications security, ACM, pp 27–37

Hoi SC, Wang J, Zhao P (2014) Libol: A library for online learning algorithms. The Journal of Machine Learning Research 15:495–499, URL <http://LIBOL.stevenhoi.org>

Lee K, Caverlee J, Webb S (2010) Uncovering social spammers: Social honeypots + machine learning. In: Proceedings of the 33rd International ACM SIGIR Conference on Research and Development in Information Retrieval, ACM, New

York, NY, USA, SIGIR '10, pp 435–442, DOI 10.1145/1835449.1835522, URL <http://doi.acm.org/10.1145/1835449.1835522>

Lee K, Eoff BD, Caverlee J (2011) Seven months with the devils: A long-term study of content polluters on twitter. In: ICWSM, Citeseer

Lee S, Kim J (2012) Warningbird: Detecting suspicious urls in twitter stream. In: NDSS

Li Y, Long PM (2002) The relaxed online maximum margin algorithm. Machine Learning 46(1-3):361–387

Matsumoto D, Hwang HS (2011) Evidence for training the ability to read microexpressions of emotion. Motivation and Emotion 35(2):181–191

Orabona F, Crammer K (2010) New adaptive algorithms for online classification. In: Advances in neural information processing systems, pp 1840–1848

Pennebaker JW, Francis ME, Booth RJ (2001) Linguistic inquiry and word count: Liwc 2001. Mahway: Lawrence Erlbaum Associates 71:2001

Rosenblatt F (1958) The perceptron: a probabilistic model for information storage and organization in the brain. Psychological review 65(6):386

Thomas K, Grier C, Ma J, Paxson V, Song D (2011a) Design and evaluation of a real-time url spam filtering

- service. In: Security and Privacy (SP), 2011 IEEE Symposium on, IEEE, pp 447–462
- Thomas K, Grier C, Song D, Paxson V (2011b) Suspended accounts in retrospect: an analysis of twitter spam. In: Proceedings of the 2011 ACM SIGCOMM conference on Internet measurement conference, ACM, pp 243–258
- Twitter (2016) Observability at twitter: technical overview, part i
- Wang D, Navathe SB, Liu L, Irani D, Tamersoy A, Pu C (2013) Click traffic analysis of short url spam on twitter. In: Collaborative Computing: Networking, Applications and Worksharing (Collaboratecom), 2013 9th International Conference Conference on, IEEE, pp 250–259
- Wang J, Zhao P, Hoi SC (2012) Exact soft confidence-weighted learning. arXiv preprint arXiv:12064612
- Webb S (2006) Introducing the webb spam corpus: Using email spam to identify web spam automatically. In: In Proceedings of the 3rd Conference on Email and AntiSpam (CEAS) (Mountain View
- Wikipedia (2016) Spamming — wikipedia, the free encyclopedia. URL <https://en.wikipedia.org/w/index.php?title=Spamming&oldid=710141595>, [Online; accessed 25-March-2016]
- Yang C, Harkreader RC, Gu G (2011) Die free or live hard? empirical evaluation and new design for fighting evolving twitter spammers. In: Recent Advances in Intrusion Detection, Springer, pp 318–337
- Yang C, Harkreader R, Zhang J, Shin S, Gu G (2012) Analyzing spammers’ social networks for fun and profit: a case study of cyber criminal ecosystem on twitter. In: Proceedings of the 21st international conference on World Wide Web, ACM, pp 71–80
- Yang L, Jin R, Ye J (2009) Online learning by ellipsoid method. In: Proceedings of the 26th Annual International Conference on Machine Learning, ACM, pp 1153–1160
- Yardi S, Romero D, Schoenebeck G, et al (2009) Detecting spam in a twitter network. *First Monday* 15(1)
- Zinkevich M (2003) Online convex programming and generalized infinitesimal gradient ascent