

Prospects for Detecting Deception on Twitter

Ulrik Franke & Magnus Rosell
 FOI – Swedish Defence Research Agency
 SE-164 90 Stockholm, Sweden
 e-mail: {ulrik.franke, magnus.rosell}@foi.se

Abstract—This paper discusses the prospects for building a system that helps a human analyst to detect deception on Twitter. First, based on a previously published taxonomy of cyber deception, Twitter deception is examined. Second, a number of indicators are introduced. These are cues, helpful for a human analyst but not necessary or sufficient characteristics of deception. Third, with the indicators as a background, the prospects for deception detection are discussed. The paper represents a first step towards a systematic method for detecting deception on Twitter.

Keywords—Twitter, deception, social network analysis, opinion mining, sentimental analysis, anomaly detection, natural language processing

I. INTRODUCTION

There is a lot of deception on the Internet. Most of it is so harmless that we hardly think of it as deception, such as the fact that users selectively post information on social networks in order to convey certain images of themselves [1]. Some of it is more annoying and can lead to poor decisions, such as hotel reviews on popular sites routinely being written by agents of the hotels themselves [2]. A small part of it is more sinister, such as the exposed cases of governments such as Russia [3] and the US [4] procuring systems designed to influence opinions on social networks by means of so called 'sock puppetry'. It is primarily the last category of deception – carried out by qualified and resourceful actors – that is the motivation for our research.

Our long-term aim is to construct a system that can aid an analyst in detecting deception on Twitter. However, as deception is a very complicated concept, we do not believe that a computer system would be able to find deception as such, but would rather act as an aid to an analyst. Although the human-machine-interaction aspects of this are indeed very important, we leave that discussion for elsewhere. Here we are mainly concerned with describing different kinds of deception and how they potentially could be detected.

The rest of the paper is organized as follows. In Section II we give examples of deception on Twitter following a general taxonomy for deception in cyberspace. We sketch a system for detecting deception in Section III. In Section IV we introduce the very wide notion of indicators that could help a human analyst to find deceptions. Section V describe

how some of the examples of deception could potentially be detected. Finally, Section VI contains some conclusions.

II. A TAXONOMY FOR TWITTER DECEPTION

Deception is a difficult subject lacking any consensus characterization. Bell & Whaley present a general deception taxonomy, based on two ways to distort reality: dissimulation, i.e. hiding the real, and simulation, i.e. showing the false [5]. In this article, we instead use the cyber-specific taxonomy proposed by Rowe [6]. It is based on linguistic case theory and offers 32 cases grouped into seven categories:

- 1) *Spatial cases* pertain to misleading locations, directions, etc.
- 2) *Time cases* pertain to misleading time stamps, records, frequencies etc.
- 3) *Participant cases* pertain to giving misleading impressions about the identity of recipients, senders, objects, beneficiaries etc.
- 4) *Causality clues* pertain to misleading causes, purposes, effects etc.
- 5) *Quality cases* pertain to misleading contents, values, measures, etc.
- 6) *Essence cases* pertain to misleading ontological features of actions, or misleading types or contexts for actions.
- 7) *Speech-act cases* pertain to misleading communication.

This taxonomy aims (i) to span the whole space of deception, and (ii) to be helpful to deception planners in brainstorming. There is no straightforward way to classify deception, and a deception action typically maps to several of the 32 cases.

Rowe also identifies *second-order deception*; i.e. deception that depends for its success on the recognition of a first-order deception. His example is an obvious denial-of-service attack masking a subtle buffer overflow attack – a decoying strategy that allegedly has been used against a number of US banks in 2012 [7].

The taxonomy concerns deception *actions*. In the case of Twitter we can split all actions into two conceptually different parts:

- The semantic action, i.e. stating the semantic content of the message.
- The posting action, i.e. the action of posting the tweet.

Someone tweeting "Obama is injured", when in fact he is not, is an example of a semantic deception. An example of a posting deception is if the tweet "Obama is injured" is widely re-tweeted making it appear trustworthy, when in fact the re-tweeters are paid. Real deceptions are often a mix of both parts.

In the following subsections we apply Rowe's taxonomy to Twitter, describing examples that often combine several of the 32 cases. The speech act cases have been left out, as Rowe himself notes that they are mostly covered by the others.

A. Space Cases

Space cases pertain to misleading locations, directions, etc. First, locations in a tweet might be misleading, and second, location metadata associated with the tweet or the tweeter might be misleading.

Direction and orientation deception. Directions mentioned in a message could be misleading. The @-notation can give a tweet a kind of direction, but this is rather an example of a recipient deception, see Section II-C.

Location-at deception. Locations in a tweet can be misleading. For instance "The suspect was spotted in Stockholm", when in fact it was Gothenburg. A tweeter can also manipulate account or tweet locations in metadata.

Location-from deception. The location of the tweeter might be considered to be from where a tweet originated. The message might contain misleading information about where something started.

Location-through deception. Tweets containing a deception can be retweeted by several tweeters, whose locations might be considered where the action has passed through. The message might contain misleading information about where something passed through.

B. Time Cases

Twitter provides a lot of detail regarding time: analogous to space, each tweet has a time-stamp. Times may also be mentioned in the tweet.

Frequency deception can be achieved using manual tweeters or bots. To make a topic seem "hotter", a bot could post tweets concerning it whenever real tweeters do not. The other time cases (*time-at deception*: time at which something occurred; *time-from deception*: time at which something started; *time-to deception*: time at which something ended; and *time-through deception*: time through which something occurred) can also be achieved using bots that change the time-frame in which a topic seems to be discussed.

It is also possible to make misleading statements in tweet messages to try to achieve any of these cases, e.g. "The army holds exercise maneuvers every week at location X".

C. Participant Cases

Agent deception on Twitter involves someone masquerading as someone else. There are least two methods: (i) *A* sets up a false account claiming to belong to *B*, or (ii) *A* obtains control (e.g. through password theft) of *B*'s actual account.

One *beneficiary deception* on Twitter is the scam: *A* crowdfunds for *B*, but the funds raised go to *C*. There is also the lie: *A* tweets that 1 000 people are gathered on square *X* supporting *B*, whereas they actually support *C*. A third variety is the conspiracy theory, where *A* tweets that event *X* (e.g. a political decision or the development of the stock market) disproportionately benefits actor or group *B* (contrary to the evidence).

A version of *experiencer deception* on Twitter is when *A* tweets a personal message to *B* and expects *B* alone to see it, but *B* also shares it with *C*, perhaps since *C* (e.g. a spouse or a colleague) just happened to be around. A more sinister version is systematic eavesdropping, when neither *A* nor *B* knows that *C* also reads the message.

Recipient deception on Twitter occurs e.g. when *A* tweets a personal message to *B*, but *C* has hacked the account of *B*. In this one-way communication example, the victim *A* to some extent fools himself. In a two-way communication example *A* tweets to *B*, believing *B* to be from nation *X*, and *B* actively responds so as to affirm this belief, though *B* is actually from nation *Y*. Thus recipient deception is not always about proper names – *B* or *C* – but sometimes about characteristics.

An example of *instrument deception* is when *A* uses an icon and a user name related to subject *X*, making *B* believe that *A* is knowledgeable on subject *X*. (Such simple methods effect perceived credibility a lot [8].)

Object deception on Twitter occurs when an action appears to be done to another object than the user believes. Some cases involve agent, beneficiary, or experiencer/recipient deception. Another case is a spam detector being fooled by different-looking shortened URLs all leading to the same web page (a classic spamming technique [9], easily adoptable to Twitter). Faked photos are also object deception: *A* tweets a photograph of protesting people on what he claims is square *A*, but it actually depicts square *B*. If the file contains false geodata, the file itself is deceptive.

D. Causality Cases

Cause deception on Twitter occurs when the cause of an action or event is misinterpreted by a user. For example, if *A* obtains control of an account belonging to *B*, thereby stopping *B* from using it, *A* could use the account in a way that makes *B* believe that the problem is technical.

An example of *contradiction deception* is when *A* writes a tweet *X* on subject *S* that contradicts what tweeter *B* usually expresses on this subject. From this, followers of *A* and *B* may assume that *A* and *B* are opponents, when in

fact they have agreed to pose as opponents on S to gain influence when they join forces on subject T .

Effect deception occurs when an effect is brought about differently than perceived by the audience. For example, if A hires well-known xenophobe B to re-tweet the tweets of C , then the general audience becomes suspicious of C through no fault of his own. The effect is achieved in a causal manner that differs from what C and the general audience believes. A more common case is when user A has thousands of followers, making B believe that A writes well on interesting topics, whereas A has actually bought the followers.

Purpose deception. A tweeter A could have many purposes for tweeting about a source X , e.g. an article. If A promotes X , promotion seems to be the purpose of A 's writing. However, if A is "trolling", the purpose might be to upset others and incite them to condemn X .

E. Quality Cases

Accompaniment and *content deception* on Twitter occurs whenever a URL purports to lead somewhere it does not. It might also be argued that the space *location-at* deception of tweeting that this is square A (when it is actually square B) and the *object* deception of a photograph with false geodata together constitute an accompaniment deception.

Manner deception on Twitter is almost indistinguishable from *effect deception* discussed above.

Material deception is rare in cyberspace, claims Rowe, because "everything is represented as bits" [6]. It might be argued, however, that some forms of object deception are close to material deception, e.g. photo manipulation.

Measure deception on Twitter occurs when numbers of followers, retweets etc. are manipulated. This can happen in several ways, such as using bots to tweet, hiring real users to tweet, placing malware in a target computer to manipulate statistics on the client side, or placing malware on a Twitter server to manipulate statistics directly.

Order deception misleads about the order in which things occur. For example, A tweets to the respected journalist B that a news story is breaking on topic X , backing this up with a deceptive URL (e.g. to a hacked news agency, to a mockup site, or to a site reporting on unrelated topic Y in a foreign language.) If B then reports on X , B might be the first reputable journalist doing so, believing himself to be the second. A similar case is when user A uses accounts A' , A'' and A''' to approach B , making it look like several independent users are retweeting a story, but B is actually the first real user to do so.

Value deception as defined by Rowe pertains to the data transmitted by the action in a software sense, i.e. arguments passed to executables etc. This is rare on Twitter, URLs with malicious software can easily be distributed. However, value deception by changing numerical values in tweets is straightforward, e.g. user A tweeting that 1 000 people are gathered on square X , whereas they are actually only 100.

F. Essence Cases

Type/supertype deception is deception in the ontological type of an action. For example, if bot A retweets user B , user B might perceive this as an act of interest, but the bot actually just did it to entice similar behavior from B , thus helping to spread the bot's message.

Whole/part deception is deception in the ontological context of an action, so that an action perceived to be part of one thing, is actually part of something else. For instance, A and B are engaged in a conversation that A believes to be normal social behavior, but B is actually gathering intelligence about A for future use.

III. A SYSTEM FOR DETECTING DECEPTION

In the remainder of this paper we will discuss possible ways to detect deception. As a rule, we do not expect a computer system to detect deception automatically – deception is by its very nature not explicit. Instead we aim to detect clues (or indicators) of deception. Our indicators are not meant to be sufficient or necessary characteristics of deception, but rather meant to be useful to the analyst. In particular, if several clues are found, they may well lead us to believe that their common cause is a deception, but such fusion of clues will most often be made by the analyst.

A system should be able to detect several different indicators and allow a user to pursue them further and/or combine them with other clues he/she instructs the system to extract. This iterative process should help the analyst sift through the huge amounts of data that even a specific part of Twitter constitutes. Also, it is important that the system allows the analyst to use background knowledge from outside the system, such as a certain user being suspicious.

IV. INDICATORS FOR TWITTER DECEPTION

Many studies describe interesting features extracted from Twitter, interesting in the context of detecting deception. In particular, three papers have inspired our selection of groups and classes of features discussed below: O'Donovan et al. [10], Castilo et al. [11], and Gupta et al. [12]. While the following is by no means an exhaustive enumeration, we try to cover several very different kinds of indicators.

A. Metadata as Indicators

Twitter supplies metadata about tweets and users, all of which is potentially interesting: Who wrote a particular tweet? When? Where does the user live? Who retweets which message by whom? etc.

B. Text Indicators

Though tweets are not "traditional" text (being very short, and informally written) natural language processing (NLP) [13] techniques are still applicable. However, they might not perform as well as on e.g. news text. Entity (e.g. *persons* and *organizations* for participant deceptions, *locations* for

space or time deception) and event extraction offer starting points for further tweet analysis.

Detecting explicit claims that someone has deceived someone else might be very useful. This borders on automatic deception/lie detection based on text features [2]. However, such work is very hard on traditional text, and probably even harder on the short tweet messages.

Information Retrieval methods [14] can also provide useful indicators concerning specific tweets or groups of tweets: keyword searches, text similarity calculations, groupings of different kinds, such as categorizations, clusterings, and topic modeling. Categorizations into opinion and sentiment classes [15] can be used for finding polarized or heated discussions.

C. Network Indicators

Social Network Analysis (SNA) [16] can be applied to Twitter using for instance (i) the follow relation and (ii) the implicit relation of who retweets who. Basic SNA gives us measures/indicators of closeness between users, user groups, and the centrality of a user in a network. In the following, we give some relevant examples. Leskovec et al. address the problem of which nodes in a network to monitor in order to quickly detect the spreading of information [17]. A similar problem is addressed by Gomez Rodriguez et al. who study how information such as certain claims or pieces of news propagate in a network [18]. Anagnostopoulos et al. have studied influence in social networks, i.e. the phenomenon when the behavior of a user induces friends to behave in a similar fashion [19]. Yamaguchi et al. propose a link-based algorithm, TURank, for finding authoritative users on Twitter based on three principles: [20]: (i) A user followed by many authoritative users is likely to be an authoritative user. (ii) A tweet retweeted by many authoritative users is likely to be a useful tweet. (iii) A user who posts many useful tweets is likely to be an authoritative user.

D. Interesting User Indicators

Some users are particularly interesting to detect (and possibly remove before analysis), e.g. users that are being paid for writing certain messages, accounts that has been hijacked, and bots. Stringhini et al. identify four different kinds of spam bots [21] and are able to detect them with very high accuracy. Chu et al. also identify *cyborgs* (bot-assisted humans or human-assisted bots) [22].

E. Polarization Indicators

Polarized discussions might be prone to exhibit deception. Guerra et al. analyze communities with respect to polarization looking specifically at the boundary nodes/users [23]. Conover et al. study tweets before the 2010 US congressional midterm elections and find big differences between the different networks of users that can be built either from mentions or from retweets [24].

F. Credibility Indicators

Credibility is closely connected to deception: Credible sources do not deceive – therefore deceivers most of all want to influence credible sources. Canini et al. try to find credible users by a hybrid approach; combining graph and text based methods [25]. Castillo et al. address the problem in greater detail in their attempt to evaluate event credibility on Twitter [11]. Gupta et al. offer a more sophisticated approach in their study of how to find credible information on current events using Twitter [12].

G. Anomaly Indicators

Looking for anomalies from what normal activity looks like by statistical means is a general strategy, applicable to any indicator. Such anomaly detection has proved successful in finding emerging topics on Twitter [26].

V. THE PROSPECTS FOR DETECTION

Based on the discussions in Sections II through III we will now briefly outline how some of the deception cases could be detected. Remember that this will be an iterative process where an analyst is aided by the computer system.

A. Space Cases

To detect the space cases, locations must be found, either from metadata or by entity recognition in the tweeted text. Obviously, locations extracted must be combined with other indicators to be indicative of deception. An example would be extracting the stated lie in "A has been lying about living in location L_1 , when in fact he lives in L_2 ." If identified and presented to an analyst, this indicates *location-at* deception.

B. Time Cases

Just as for space cases, times must be extracted either from metadata and/or by entity recognition in the tweeted text. Again, times extracted must be combined with other indicators to be indicative of deception, e.g. a time stated in a tweet message by a suspicious tweeter might be misleading.

C. Participant Cases

Detecting deception regarding participants first requires detection of those participants (e.g. persons, organizations, countries etc.) either from metadata or from the tweeted text as in the previous sections.

Consider *agent deception* where A controls B 's account. If the take-over is accompanied by a change in the temporal activity profile of the hijacked account [27] or if B suddenly doubles the number of people he follows, this can be detected. Such an anomaly indicator would be particularly powerful if combined with a black-list of suspicious users.

Consider the *beneficiary deception* where A tweets that 1 000 people are gathered on square X in support of B , while the people actually support C . If others also tweet about the people on square X , the discrepancy might be

detected through NLP analysis of the semantics of the tweets.

The case of *object deception* with faked or re-used photos can be detected by a system that intercepts tweets and tests images using a service such as Google image search, flagging re-used photos as suspicious. Geodata differing from other versions of the same photo can also be flagged.

D. Causality Cases

Consider the *effect deception* where user *A* has thousands of, mostly bought, followers. A graph analysis of the followers can expose this – if many of them appear to be following very different users, it might be suspected that they do not follow out of interest, but rather for profit. The better the normal follower-pattern is known, the better this indicator becomes.

E. Quality Cases

Accompaniment and *content deception* with shortened URLs can be exposed by a system that intercepts tweets and follows URLs in a sandboxed environment. URLs that look different but point to the same resource can be flagged.

Consider the *order deception* when journalist *B* is being approached by *A'*, *A''* and *A'''* and breaks a story that believes is already in the news. A system that investigates accounts and flags those that are newly created can expose the simplest case. A more sophisticated attack with more reputable fake users "bred" for sale is more challenging, though a graph analysis might still show discrepancies from typical patterns exhibited by real users.

F. Essence Cases

Type/super type deception as described in Section II-F is possible to detect if one is able to find bots. Any conversation by a bot is a potential deception. An analyst may determine the message a bot is trying to spread if provided with (a summary of) the conversations a bot has had.

Whole/part deception as described in Section II-F is harder to detect. A system could provide an analyst with the conversations an interesting tweeter *A* has had. If *A* is a bot, trawling other users for personal information, a pattern in the conversations might become apparent to the analyst given enough data.

VI. CONCLUSIONS

This paper has explored the prospects for building a system that helps a human analyst to detect deception on Twitter. The contribution has several stages: First, using a previously published taxonomy of cyber deception, we examined how Twitter can be used as a vector for deception. Though the cases identified are probably not exhaustive, it is reasonable to believe that large and relevant portions of potential deception strategies were found. Second, we introduced the wide notion of indicators. These are clues that

can be helpful for a human analyst and are not necessary or sufficient characteristics of deception. Third, with the indicators as a background, the prospects for deception detection were discussed. Though the analysis needs empirical verification, it is based on structured and literature-based reasoning about what to find and how to find it.

In conclusion, we believe that a system could aid a human analyst in the detection of deception if he or she is provided with the right indicators and allowed to interact with the system in an iterative manner.

ACKNOWLEDGMENT

This work was supported by the R&D program of the Swedish Armed Forces. Christian Mårtensson and Fredrik Johansson gave valuable comments on the article.

REFERENCES

- [1] J. Peluchette and K. Karl, "Examining students intended image on facebook:what were they thinking?!", *Journal of Education for Business*, vol. 85, no. 1, pp. 30–37, 2009.
- [2] M. Ott, Y. Choi, C. Cardie, and J. T. Hancock, "Finding deceptive opinion spam by any stretch of the imagination," in *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies-Volume 1*. Association for Computational Linguistics, 2011, pp. 309–319.
- [3] I. Barabanov, I. Safronov, and E. Chernenko, "Razvedka botom [Intelligence Using a Bot]," <http://www.kommersant.ru/doc/2009256>, Aug. 2012, Kommersant, retrieved 7 November 2012.
- [4] N. Fielding and I. Cobain, "Revealed: US spy operation that manipulates social media," <http://www.guardian.co.uk/technology/2011/mar/17/us-spy-operation-social-networks>, Mar. 2011, the Guardian, retrieved 18 January 2013.
- [5] J. B. Bell and B. Whaley, *Cheating and deception*. Transaction Publishers New Brunswick, 1991.
- [6] N. Rowe, "A taxonomy of deception in cyberspace," in *International Conference on Information Warfare and Security*, 2006, pp. 173–181.
- [7] J. Robertson, "Cheapest Way to Rob Bank Seen in Cyber Attack Like Hustle," <http://www.bloomberg.com/news/2013-05-06/cheapest-way-to-rob-bank-seen-in-cyber-attack-like-hustle.html>, bloomberg, published May 6, 2013, retrieved November 29 2013.
- [8] M. R. Morris, S. Counts, A. Roseway, A. Hoff, and J. Schwarz, "Tweeting is believing?: understanding microblog credibility perceptions," in *Proceedings of the ACM 2012 conference on Computer Supported Cooperative Work*. ACM, 2012, pp. 441–450.
- [9] J. Ratkiewicz, M. Conover, M. Meiss, B. Gonçalves, A. Flammini, and F. Menczer, "Detecting and tracking political abuse in social media," in *Proc. 5th International AAAI Conference on Weblogs and Social Media (ICWSM)*, 2011.

- [10] J. O'Donovan, B. Kang, G. Meyer, T. Höllerer, and S. Adalii, "Credibility in context: An analysis of feature distributions in twitter." in *SocialCom/PASSAT*. IEEE, 2012, pp. 293–301.
- [11] C. Castillo, M. Mendoza, and B. Poblete, "Information credibility on twitter," in *Proceedings of the 20th international conference on World wide web*. ACM, 2011, pp. 675–684.
- [12] M. Gupta, P. Zhao, and J. Han, "Evaluating event credibility on Twitter," in *2012 SIAM International Conference on Data Mining*. SIAM, 2012, pp. 153–164.
- [13] D. Jurafsky and J. H. Martin, *Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics and Speech Recognition*, andra ed., ser. Prentice Hall series in artificial intelligence. Upper Saddle River, New Jersey: Prentice Hall, Pearson Education International, 2009.
- [14] C. D. Manning, P. Raghavan, and H. Schütze, *Introduction to Information Retrieval*. New York: Cambridge University Press, 2008.
- [15] B. Pang and L. Lee, "Opinion mining and sentiment analysis," *Foundations and Trends in Information Retrieval*, vol. 2, no. 1–2, pp. 1–135, Jan. 2008.
- [16] S. Wasserman and K. Faust, *Social Network Analysis: methods and applications*. Cambridge University Press, 1994.
- [17] J. Leskovec, A. Krause, C. Guestrin, C. Faloutsos, J. VanBriesen, and N. Glance, "Cost-effective outbreak detection in networks," in *Proceedings of the 13th ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, 2007, pp. 420–429.
- [18] M. Gomez Rodriguez, J. Leskovec, and A. Krause, "Inferring networks of diffusion and influence," in *Proceedings of the 16th ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, 2010, pp. 1019–1028.
- [19] A. Anagnostopoulos, R. Kumar, and M. Mahdian, "Influence and correlation in social networks," in *Proceedings of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, 2008, pp. 7–15.
- [20] Y. Yamaguchi, T. Takahashi, T. Amagasa, and H. Kitagawa, "TURank: Twitter user ranking based on user-tweet graph analysis," in *Web Information Systems Engineering–WISE 2010*. Springer, 2010, pp. 240–253.
- [21] G. Stringhini, C. Kruegel, and G. Vigna, "Detecting spammers on social networks," in *Proceedings of the 26th Annual Computer Security Applications Conference*, ser. ACSAC '10. New York, NY, USA: ACM, 2010, pp. 1–9.
- [22] Z. Chu, S. Gianvecchio, H. Wang, and S. Jajodia, "Who is tweeting on twitter: Human, bot, or cyborg?" in *Proceedings of the 26th Annual Computer Security Applications Conference*, ser. ACSAC '10. New York, NY, USA: ACM, 2010, pp. 21–30.
- [23] P. H. C. Guerra, W. M. Jr., C. Cardie, and R. Kleinberg, "A measure of polarization on social media networks based on community boundaries." in *ICWSM*, E. Kiciman, N. B. Ellison, B. Hogan, P. Resnick, and I. Soboroff, Eds. The AAAI Press, 2013.
- [24] M. D. Conover, J. Ratkiewicz, M. Francisco, B. Gonçalves, A. Flammini, and F. Menczer, "Political polarization on twitter," in *Proc. 5th Intl. Conference on Weblogs and Social Media*, 2011.
- [25] K. R. Canini, B. Suh, and P. L. Pirolli, "Finding credible information sources in social networks based on content and social structure," in *Privacy, security, risk and trust (PASSAT), 2011 IEEE third international conference on and 2011 IEEE third international conference on social computing (SOCIALCOM)*. IEEE, 2011, pp. 1–8.
- [26] T. Takahashi, R. Tomioka, and K. Yamanishi, "Discovering emerging topics in social streams via link anomaly detection," in *Data Mining (ICDM), 2011 IEEE 11th International Conference on*. IEEE, 2011, pp. 1230–1235.
- [27] F. Johansson, L. Kaati, and A. Shrestha, "Detecting multiple aliases in social media," in *Proceedings of the 2013 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining*. ACM, 2013, pp. 1004–1011.