

Industrial Uses for Authorship Analysis

Patrick Juola

Evaluating Variations in Language Laboratory

Duquesne University

600 Forbes Avenue

Pittsburgh, PA 15282 USA

Email: juola@mathcs.duq.edu

Tel: +1 (412) 396-2276

Abstract—Text classification is an important technology to help industry handle millions of words of customer communications. We discuss a less well-known application of text technology, specifically authorship attribution and profiling. By examining not the content but the style of these communications, computers can learn not only what people are writing about, but things about the people writing as well, such as their identity, demographics, and even psychometrics. We provide several applications to illustrate the value of this important emerging technology.

I. INTRODUCTION

Sometimes knowing *who* wrote a document is as important as knowing *what* was written. Sometimes you need to know who your critics are, even when they hide behind anonymous sounding names on Internet forums. Sometimes you need to know who actually sent a letter, a piece of email, or a text message. While the idea of using computers to analyze the contents of a document is well known, the idea of using one to analyze the author is perhaps less well-understood. In this paper, we provide examples both of how this kind of analysis can be done, and more importantly, of why this capacity is important.

II. BACKGROUND

A. Text Categorization

More text information is available now than ever before in history, and no human can possibly read it all. At the same time, through channels like Internet forums and product commentary, customers have more ability to influence each other than ever before, and suppliers need to be able to read the same documents. But if not humans, then...computers? By reading documents at the speed of electronics, emerging technologies make it possible to keep up with the flood of text.

Two key examples of this technology are topic modeling and sentiment analysis. Topic modeling [1] infers “the main themes that pervade a large and otherwise unstructured collection of documents,” in order to “organize the collection according to the discovered themes.” Sentiment analysis [2] categorizes documents by the emotions and/or opinions expressed in it, usually based on a positive/negative “polarity,” describing whether the author likes/approves of the topic. For example, “awesome” is generally a positive term, as is (less obviously) “zest.” “Abdicate” is generally negative. Using these technologies together, a computer can analyze thousands or millions of comments about a product and tell the company

both *what* aspects of the product people are discussing as well as *how* they feel about the various aspects.

What this technology will not tell you is *who* is writing these comments. At one level, this is simply useful marketing information; if everyone who likes your product are college-age females, this suggests both new markets that might be opened as well as efficient channels to build the existing market. More subtly, it’s hard to tell how many actual people are behind the comments and whether the criticisms are a genuine issue or simply one hard-to-satisfy customer with too much free time. Indeed, review spam [3], [4] where praise and criticisms are written and published on a commercial scale, is a problem of increasing urgency. Another application of text categorization, stylometry (also called stylometrics) provides a potential solution to these problems. By examining the individual style of the individual writers, the computer can tell you not only about the contents of the documents, but also things about the person who wrote it.

B. Theory of Stylometry

So how does this work? The basic theory of traditional stylistics is fairly simple. As McMenamin describes it,

At any given moment, a writer picks and chooses just those elements of language that will best communicate what he/she wants to say. The writer’s “choice” of available alternate forms is often determined by external conditions and then becomes the unconscious result of habitually using one form instead of another. Individuality in writing style results from a given writer’s own unique set of habitual linguistic choices.[5]

Coulthard’s description is similar:

The underlying linguistic theory is that all speaker/writers of a given language have their own personal form of that language, technically labeled an idiolect. A speaker/writers idiolect will manifest itself in distinctive and cumulatively unique rule-governed choices for encoding meaning linguistically in the written and spoken communications they produce. For example, in the case of vocabulary, every speaker/writer has a very large learned and stored set of words built up over many years. Such sets may differ slightly or considerably from the word sets that all other speaker/writers have similarly

built up, in terms both of stored individual items in their passive vocabulary and, more importantly, in terms of their preferences for selecting and then combining these individual items in the production of texts. [6]

A non-obvious but key application is to the legal system. For example, a famous dispute over the ownership of a significant part of Facebook (*Ceglia v. Zuckerberg and Facebook*) depended in part upon a set of disputed writings. These writings were email, allegedly written by Mark Zuckerberg, that purported to show that Paul Ceglia owned half of Facebook. Of course, if these writings were not by Zuckerberg, they showed nothing of the sort. McMenamin's report analyzed eleven different and distinct "features" of the writing in both the known (undisputed) email and the disputed email. One feature, for example, hinged on the spelling of the word *cannot*, and in particular whether it was written as one word (*cannot*) or as two (*can not*). Another feature was the use of the single word "sorry" as a sentence opener (as opposed, for example, to "I'm sorry"). [5] submitted a report that showed that the writing style of a set of undisputed email (that Zuckerberg acknowledged having written) differed in a number of important ways from the disputed writings, and concluded that "[i]t is probable that Mr. Zuckerberg is not the author of the QUESTIONED writings." (Capitalization in original.)

Similarly, [6] describes a (redacted) case of authorship of a disputed email leaked from a company under questionable circumstances. Coulthard similarly discussed (among other features) the use of the specific phrase "disgruntled employees." [7] describes a case of potential murder, where the authorship of a set of SMS (text) messages found on a cell phone constituted a key element in establishing both the time of death (when the writing style of these messages shifted radically) and showed strong indications of an attempt to cover up the murder via arson. By examining features including variant spellings such as "wiv" for "with" and "wud" for "would," he was able to show key differences between the writing of the messages and the typical writing of the phone's owner. He was also able to show key similarities between the writings of the (alleged) murderer/arsonist and one of the suspects in the case.

In other examples, [8] describes a case in immigration court, where an applicant for political asylum was able to lay claim to a number of anonymous newspaper columns critical of his home government, and therefore establish a reasonable fear of persecution upon return to his homeland. [9] describes another murder case, one where the crime scene included a suicide note typed on a shared computer, but stylistic analysis was able not only to show that it had not been written by the decedent, but also to identify someone else as the killer.

Computer-based stylometry applies the same general theory, but with a few major differences. The basic assumption that people make individual choices about language still holds, but instead of *ad hoc* features selected by examination of the specific documents, the analysts use more general feature sets that apply across the spectrum of problems. One common feature set is the frequency of common words such as articles and prepositions [10], [11], [12]. Because these words tend

both to be common and also not to carry strong semantic associations, their frequencies tend to be stable across documents and genres, but these frequencies can also be shown to vary strongly across individuals. Another commonly used feature set is the frequency of common groups of consecutive words (word *n*-grams) or consecutive characters (character *n*-grams) [13], [14], [15]. Using these feature sets or others [16], the features present in a document are automatically identified, gathered into collections of feature representations (such as vector spaces), and then classified using ordinary machine learning algorithms [17], [18], [19] to establish the most likely author.

A particularly good example is Binongo's study of the *Oz* books [11]. The backstory is fairly simple: the series was started with L. Frank Baum's publication of *The Wonderful Wizard of Oz* and continued until his death in 1919. After his death, the publishers asked Ruth Plumly Thompson to finish "notes and a fragmentary draft" of what would become *The Royal Book of Oz*, the 15th in the series, and then Thompson herself continued the series until 1939, writing nearly twenty more books. The underlying question is the degree to which this "fragmentary draft" influenced Thompson's writing; indeed, scholars have no evidence that the draft ever existed. Binongo collected frequency statistics on the fifty most frequent function words across the undisputed samples and analyzed them using principal component analysis (PCA). Reducing these fifty variables down to their first two principal components produced an easily graphable distribution that showed clear visual separation between the two authors. When the *Royal Book* was plotted on the same scale, it was shown clearly to lie on Thompson's side of the graph, confirming that "from a statistical standpoint, [the *Royal Book*] is more likely to have been written in Thompson's hand."

III. APPLICATIONS

A. Attribution

In 1996, the novel *Primary Colors* was published. A roman-à-clé purporting to describe Clinton's 1992 presidential campaign, it provided an insightful view into late 20th century American politics. Or did it? If the anonymous author actually had inside knowledge, that was one thing. On the other hand, if it was just a potboiler by an ordinary novelist, it may no more accurately have reflected reality than a Spider-man comic book describes life in contemporary New York City. As part of the discussion surrounding this book, linguist Don Foster [20] showed that the writing was very similar to that of columnist Joe Klein, who later acknowledged authorship. Another recent high-profile example [21], [22] is that of the author of *A Cuckoo's Calling*, by Robert Galbraith. Although Galbraith was a first-time author, numerous critics noted that the authorial voice was unusually polished and confident. Formal analyses of writing style, performed at the behest of the *Sunday Times*, later identified [23] J.K. Rowling, author of the *Harry Potter* books, as the actual author. Literature scholars have been interested in questions of authorship for centuries, as typified by the discussions of authorship of Biblical book of Acts [24], traditionally ascribed to the author of the book of Luke, and of the authorship of the *Illiad* and the *Odyssey*, still an open question [25]. However, identifying the author of a document can be of interest to other parties as well.

Another common application is journalism. As with the Rowling case [21], many questions arise from a matter of public interest, driven by journalists. Another recent example is *Newsweek's* analysis of the Bitcoin design documents, officially written by a person named “Satoshi Nakamoto” (which may have been a pseudonym), and attributed by *Newsweek* to a retired engineer named Dorian Nakamoto. Stylometric analysis of these documents [26] against an appropriate set of known documents showed “that Dorian Nakamoto was not found to be a plausible candidate author, and in fact, one of the distractor authors (Neal J. King) was found to be a better match to Satoshi Nakamoto than any other distractor or than Dorian.” [27]

B. Profiling

A related problem is that of authorship profiling [28], the study of other authorial characteristics such as gender, age, education level, native language, personality and so forth. [29], [30], [31], [32], [33], [34], [35], [36], [37], [38], [39] Profiling is in some ways a more important problem than attribution. Profiling can be and is used [28] on a larger scale to infer group properties of a large number of people.

Profiling is done in the same way as attribution, but instead of offering training documents labeled by author, the system is provided with documents representing specific groups, such as essays written by college graduates and by non-college graduates, or by speakers of UK and US English. The same feature selection and classification techniques will infer the appropriate markers for group membership and classify novel documents accordingly. (To illustrate, an obvious feature for distinguishing UK vs. US dialects would be vocabulary, and specifically items like “lorry,” “ironmonger,” and “tarmac.”)

Authorship profiling has obvious commercial potential (what can I learn about the people who post negative reviews of my product?) but is also of significant interest to other fields, such as law enforcement. Among other applications, it forms one of the technologies underlying DARPA’s Active Authentication project [40], [15], [41], based on the theory that if I write (or more generally, interact with the computer) like an introvert, but the person actually at the keyboard behaves as an extrovert, then that person is probably not me. In the event of an actual security incident, learning *about* the intruder can provide a useful start for investigation and response. Other applications may include telemedicine, for example, allowing the nonintrusive identification and assessment of risk factors such as bipolar disorder [42], low self-esteem [43], depression [44] or suicidality [45].

The methodological basis of these analyses are very similar to the authorial analysis, and the same software can be used for both applications [32], [46]. Indeed, in many cases [15] very similar feature spaces and classification methods are among the best-performing; the only difference is in the labeling of the training corpus. There are several proofs-of-concept in this space [32], [47], [38], [43], [42], illustrating that it is quite practical to do this kind of profiling for a number of different attributes, including both normal [48], [29], [47], [38], [43] and pathological [44], [45], [42] psychological traits as well as ordinary demographic information [32], and even handedness [15]. This technology has even been used to infer deceptive intentions [49], [50], [51].

IV. THREE NOVEL APPLICATIONS

A. Case 1: Identifying commercial sock puppets

We have identified in the previous paragraphs all of the necessary ingredients to begin addressing a key problem of the modern commercial world, that of commercial deceptive social media. Deceptive customer reviews—whether paid-for-positive reviews by shills, or damning reviews placed by agents of the competition—are becoming a major issue in e-commerce and a major problem for businesses whose primary product is review aggregation. Deceptive review spam is used as a marketing tool by corporations [3], [52], political pressure groups [53], and even national governments [4], but can also simply be the acts of a single active person with an axe to grind.

As discussed in the previous section, this deceptive intention can be detected, as can posts by the same author using multiple identifiers and user names. This provides a relatively simple way to allow an analyst to disregard multiple postings or deceptive postings. In fact, it would even be straightforward for an aggregator such as Yelp to eliminate these from consideration in offering “average” customer ratings, or for law enforcement such as the Federal Trade Commission to initiate proceedings as appropriate. By identifying overrepresented (or outright deceptive) comments, this enables the merchant to develop a more representative picture of the customer base and take actions grounded in a better and more realistic understanding of the true situation.

B. Case 2: Identity as a behavioral biometric

Passwords are generally considered to provide weak security. [40] They can be forgotten, guessed, or stolen. More subtly, passwords only provide momentary security up-front, at login time. When the user gets up to get coffee, the computer retains the user’s credentials and will continue to provide access to anyone who sits down at the keyboard. Chaski [9] provides an example of a legal dispute hinging on who actually sent inappropriate email from a (shared) corporate computer, but a more common problem might involve insider threats, where someone uses someone else’s leftover credentials to access beyond his/her authorization.

In 2012, DARPA [40] proposed to develop “Active Authentication,” an alternative approach to computer security where users are continuously and actively reassessed on an ongoing basis to determine whether or not they are still authorized to use the computer. In the event that a user does something that causes the system to question their identity, a security alert can be raised (and appropriate action taken).

One of the technologies assessed for this project is authorship attribution and profiling as a form of linguistic biometric [15], [41]. In simple terms, if the person writing email is not writing the way the authorized user would write, then the person writing may not be the authorized user. Similarly, if a person is drafting a document in the wrong writing style, there may be an issue. Even profiling can be applied to this issue — if the person writing the document writes like a member of the wrong group, a group to which the authorized user does not belong (and which the person normally does not write like), there may again be an issue.

Juola et al. [15] have shown this to be feasible. This group collected information about writing style from a group of 80 participants in a simulated work environment. Each person, over a one-week period, was asked to do a long-term blogging task intermixed with smaller, more explicitly-defined writing tasks of a few hours each. Using ordinary stylometric technology, they were able to identify specific participants by their writing style with roughly 60% accuracy based on as few as 500 characters, and to identify personality categories such as introversion/extraversion with approximately 80% accuracy. (By contrast, the chance baseline for identifying specific people is approximately 1%, one in 80, and for identifying personality traits approximately 50%.)

C. Case 3: Psychometrically-informed advertising and customer relations

The idea of using authorship profiling to identify demographic information about actual and potential customers has been discussed in a previous section, and using this technology to infer demographics is well-understood [54], [28]. However, demographics is only half the story, and as the active authentication project has shown, it's possible to infer mental and social traits as well as demographic ones. This makes it both possible and practical to narrowcast messages to specific people based on previous text interactions with them.

The idea of targeting advertisements (or other corporate communications) to a specific person is of course not new; that's one of the fundamental premises behind cookie-based marketing. However, author profiling technology creates new opportunities for analysis, with a new channel providing classification information without needing to gather data from external sources. One specific application for this is in "inside sales," where communications with existing clients can be (re)analyzed to determine both the best approach for maintaining and extending the relationship. This can help, for example, by allowing better matching of successful sales representatives to customers, based on the types of customers and the types of representatives. This assures that customers have representatives that will be able to connect well with them, understand their needs, and create a closer, more beneficial association. This could be done on the basis of demographic and personality data as described in the previous section, or even on the basis of ad-hoc categories for each representative, representing empirically what each representative's strengths and weaknesses are.

V. CONCLUSIONS

Text analysis is well understood as a key business technology; it lets companies deal with large sets of documents easily and efficiently. Authorship analysis is not as well known or understood, but provides another key capacity, the ability to deal with large sets of clients and customers easily and efficiently. In this paper, we have described the basics of this technology and outlined several specific, practical applications that can have a major effect on industry.

REFERENCES

- [1] D. M. Blei, "Probabilistic topic models," *Communications of the ACM*, vol. 55, no. 4, pp. 77–84, April 2012.
- [2] J. Wiebe, T. Wilson, and C. Cardie, "Annotating expressions of opinions and emotions in language," *Language Resources and Evaluation*, vol. 39, pp. 165–210, 2005.
- [3] P. Elmer-DeWitt, "Samsung fined \$340,000 for astrourfing in Taiwan," *Fortune*, 2013.
- [4] B. Sterling, "The Chinese online 'water army'," *Wired*, vol. June, 2010.
- [5] G. McMenamin, "Declaration of Gerald McMenamin," Available online at <http://www.scribd.com/doc/67951469/Expert-Report-Gerald-McMenamin>, 2011.
- [6] M. Coulthard, "On admissible linguistic evidence," *Journal of Law and Policy*, vol. XXI, no. 2, pp. 441–466, 2013.
- [7] T. Grant, "Txt 4n6: Describing and measuring consistency and distinctiveness in the analysis of SMS text messages," *Journal of Law and Policy*, vol. XXI, no. 2, pp. 467–494, 2013.
- [8] P. Juola, "Stylometry and immigration: A case study," *Journal of Law and Policy*, vol. XXI, no. 2, pp. 287–298, 2013.
- [9] C. E. Chaski, "Who's at the keyboard: Authorship attribution in digital evidence investigations," *International Journal of Digital Evidence*, vol. 4, no. 1, p. n/a, 2005, electronic-only journal: <http://www.ijde.org>, accessed 5.31.2007.
- [10] J. F. Burrows, "'an ocean where each kind...' : Statistical analysis and some major determinants of literary style," *Computers and the Humanities*, vol. 23, no. 4-5, pp. 309–21, 1989.
- [11] J. N. G. Binongo, "Who wrote the 15th book of Oz? an application of multivariate analysis to authorship attribution," *Chance*, vol. 16, no. 2, pp. 9–17, 2003.
- [12] D. L. Hoover, "Delta prime?" *Literary and Linguistic Computing*, vol. 19, no. 4, pp. 477–495, 2004.
- [13] E. Stamatatos, "On the robustness of authorship attribution based on character n-gram features," *Journal of Law and Policy*, vol. XXI, no. 2, pp. 420–440, 2013.
- [14] G. K. Mikros and K. Perifanos, *Authorship attribution in Greek tweets using multilevel author's n-gram profiles*. Palo Alto, California: AAAI Press, 2013, pp. 17–23. [Online]. Available: <http://www.aaai.org/ocs/index.php/SSS/SSS13/paper/view/5714/5914>
- [15] P. Juola, J. I. Noecker Jr, A. Stolerman, M. V. Ryan, P. Brennan, and R. Greenstadt, "Keyboard behavior-based authentication for security," *IT Professional*, vol. 15, pp. 8–11, 2013.
- [16] J. Rudman, "On determining a valid text for non-traditional authorship attribution studies : Editing, unediting, and de-editing," in *Proc. 2003 Joint International Conference of the Association for Computers and the Humanities and the Association for Literary and Linguistic Computing (ACH/ALLC 2003)*, Athens, GA, May 2003.
- [17] T. Joachims, *Learning to Classify Text Using Support Vector Machines*. Kluwer, 2002.
- [18] J. R. Quinlan, *C4.5: Programs for Machine Learning*. Morgan Kaufman, 1993.
- [19] M. L. Jockers and D. Witten, "A comparative study of machine learning methods for authorship attribution," *LLC*, vol. 25, no. 2, pp. 215–23, 2010.
- [20] D. Foster, *Author Unknown: Tales of a Literary Detective*. New York: Holt Paperbacks, 2001.
- [21] R. Brooks and C. Flynn, "JK Rowling: The cuckoo in crime novel nest," *Sunday Times*, vol. 14 July, 2013.
- [22] R. Brooks, "Whodunnit? JK Rowling's secret life as wizard crime writer revealed," *Sunday Times*, vol. 14 July, 2013.
- [23] P. Juola, "How a computer program helped reveal J. K. Rowling as author of A Cuckoo's Calling," *Scientific American*, vol. August, 2013. [Online]. Available: <http://www.scientificamerican.com/article/how-a-computer-program-helped-show-jk-rowling-write-a-cuckoos-calling/>
- [24] C. A. Evans, *Luke*, ser. Understanding the Bible Commentary Series. Grand Rapids, Michigan: Baker Books, 1990.
- [25] R. Garland, *Ancient Greece: Everyday Life in the Birthplace of Western Civilization*. New York: Sterling, 2008.
- [26] M. Herper, "Linguist analysis says Newsweek named the wrong man as Bitcoin's creator," *Forbes Magazine*, vol. March 10, 2014. [Online]. Available: <http://www.forbes.com/sites/matthewherper/2014/03/10/data-analysis-says-newsweek-named-the-wrong-man-as-bitcoins-creator/>

- [27] P. Juola, "The Rowling case: A proposed standard protocol for authorship attribution," in *Proceedings of Digital Humanities 2014*, Lausanne, Switzerland, 2014.
- [28] S. Argamon, M. Koppel, J. W. Pennebaker, and J. Schler, "Automatically profiling the author of an anonymous text," *Communications of the ACM*, vol. 52, no. 2, pp. 119–123, February 2009.
- [29] S. Argamon, S. Dhawle, M. Koppel, and J. W. Pennebaker, "Lexical predictors of personality type," in *Proceedings of the Classification Society of North America Annual Meeting*, 2005. [Online]. Available: citeseer.ist.psu.edu/744868.html
- [30] R. H. Baayen, H. van Halteren, A. Neijt, and F. Tweedie, "An experiment in authorship attribution," in *Proceedings of JADT 2002*. St. Malo: Université de Rennes, 2002, pp. 29–37.
- [31] M. Corney, O. de Vel, A. Anderson, and G. Mohay, "Gender-preferential text mining of e-mail discourse," in *Proceedings of Computer Security Applications Conference*, 2002, 2002, pp. 282–289.
- [32] P. Juola and H. Baayen, "A controlled-corpus experiment in authorship attribution by cross-entropy," *Literary and Linguistic Computing*, vol. 20, pp. 59–67, 2005.
- [33] M. Koppel, S. Argamon, and A. R. Shimoni, "Automatically categorizing written texts by author gender," *Literary and Linguistic Computing*, vol. 17, no. 4, pp. 401–412, 2002, doi:10.1093/lc/17.4.401.
- [34] T. Kucukyilmaz, B. B. Cambazoglu, C. Aykanat, and F. Can, "Chat mining for gender prediction," *Lecture Notes in Computer Science*, vol. 4243, p. 274283, 2006.
- [35] M. Oakes, "Text categorization: Automatic discrimination between US and UK English using the chi-square text and high ratio pairs," *Research in Language*, vol. 1, pp. 143–156, 2003.
- [36] H. van Halteren, R. H. Baayen, F. Tweedie, M. Haverkort, and A. Neijt, "New machine learning methods demonstrate the existence of a human stylome," *Journal of Quantitative Linguistics*, vol. 12, no. 1, pp. 65–77, 2005.
- [37] B. Yu, Q. Mei, and C. Zhai, "English usage comparison between native and non-native English speakers in academic writing," in *Proceedings of ACH/ALLC 2005*, Victoria, BC, Canada, 2005.
- [38] J. Noecker Jr, M. Ryan, and P. Juola, "Psychological profiling through textual analysis," *LLC*, vol. 28, no. 3, pp. 382–387, 2013.
- [39] G. K. Mikros, *Authorship Attribution and Gender Identification in Greek Blogs*. Belgrade: Academic Mind, 2013, pp. 21–32.
- [40] R. P. Guidorizzi, "Security: Active Authentication," *IT Professional*, vol. 15, pp. 4–7, 2013.
- [41] S. Acharya, A. Fridman, P. Brennan, P. Juola, and R. Greenstadt, "User authentication through biometric sensors and decision fusion," in *47th Annual Conference on Information Sciences and Systems (CISS 2013)*, 2013.
- [42] J. I. Noecker Jr. and P. Juola, "Stylometric identification of manic-depressive illness," in *Proceedings of DHCS 2014*, 2014.
- [43] P. Juola and J. I. Noecker Jr., "Inferring self-esteem from keyboard behavior," in *Proceedings of DHCS 2014*, 2014.
- [44] S. Rude, E. Gortner, and J. Pennebaker, "Language use of depressed and depression-vulnerable college students," *Cognition and Emotion*, vol. 18, pp. 1121–1133, 2004.
- [45] C. Poulin, B. Shiner, P. Thompson, L. Vepstas, Y. Young-Xu, B. Goertzel, B. Watts, L. Flashman, and T. McAllister, "Predicting the risk of suicide by analyzing the text of clinical notes," *PLoS One*, vol. 9, no. 3, 2014.
- [46] P. Juola, "20,000 ways not to do authorship attribution and a few that work," in *Proceedings of 2009 Biennial Conference of the International Association of Forensic Linguists (IAFL-09)*, Amsterdam, 2009.
- [47] C. Gray and P. Juola, "Personality identification through on-line text analysis," in *Proceedings of the 2011 Chicago Colloquium on Digital Humanities and Computer Science*, Chicago, IL, 2011.
- [48] J. W. Pennebaker and L. A. King, "Linguistic styles: Language use as an individual difference," *Journal of Personality and Social Psychology*, vol. 77, pp. 1296–1312, 1999.
- [49] M. Newman, J. Pennebaker, D. Berry, and J. Richards, "Lying words: Predicting deception from linguistic style," *Personality and Social Psychology Bulletin*, vol. 29, pp. 665–675, 2003.
- [50] I. Picornell, "The flexible liar – a new approach to detecting deception in written narratives," in *Forensic Linguistics: Bridging the Gap(s) Between Language and the Law: 2012 Regional Meeting of the IAFL*. Porto, Portugal: International Association of Forensic Linguists, 2012.
- [51] J. I. Noecker Jr and P. Juola, "Spanish lies : A computational reanalysis of a spanish corpus for detecting deception," in *Proceedings of the 11th Biennial Conference on Forensic Linguistics/Language and Law of the International Association of Forensic Linguists (IAFL 2013)*, Mexico City, MX, June 2013.
- [52] D. Folkenflik, *Murdoch's World*. PublicAffairs, 2013.
- [53] K. Kleiner, "Bogus grass-roots politics on Twitter," *Technology Review*, 2010.
- [54] H. van Halteren, "Author verification by linguistic profiling: An exploration of the parameter space," *ACM Transactions on Speech and Language Processing*, vol. 4, p. n/a, 2007.