

SCUOLA DI SCIENZE

Corso di Laurea Magistrale in Informatica

On
Authorship
Attribution

Relatore:
Chiar.mo Prof.
DANILO MONTESI

Presentata da:
Gabriele Calarota

Correlatore:
Dott.
FLAVIO BERTINI

Sessione III
Anno Accademico 2019-2020

*“When I was in college,
I wanted to be involved in things that would change the world”*
Elon Musk

SOMMARIO

ABSTRACT

CONTENTS

1	Introduction	11
1.0.1	Motivation and Problem Statement	11
1.0.2	Thesis Structure	11
2	Authorship attribution's tasks	13
2.1	History of methodologies	14
2.2	Training's approach	15
2.2.1	Profile-based approach	16
2.2.2	Instance-based approach	17
2.3	The real problem	18
2.3.1	Profiling problem	19
2.3.2	Needle-in-hay-stack problem	19
2.3.3	Verification problem	20
2.4	Identify the problem	20
2.4.1	Single domain vs Cross domain	20
2.4.2	Closed set vs Open set	21
3	Text characteristics analysis	23
3.1	Character Features	24
3.1.1	Affix n-grams	24
3.1.2	Word n-grams	25
3.1.3	Punctuation n-grams	25
3.2	Lexical Features	26
3.2.1	Bag of Words	26
3.2.2	Word N-grams	29
3.2.3	Vocabulary Richness	31
3.2.4	Stylometric features	32
3.2.5	Function Words	33
3.2.6	Tf-Idf	34
3.3	Syntactic Features	35

3.4	Semantic Features	36
3.4.1	Positivity and Negativity index	36
3.5	Application Specific Features	38
3.5.1	Vector embeddings of words (Word2Vec)	38
3.5.1.1	Skip-gram	39
3.5.1.2	CBOW	39
3.5.2	Vector embeddings of documents (Doc2Vec)	40
4	State of the art: Data collection and Techniques for Authorship Attribution	43
4.1	SVM	44
4.1.1	SVM for authorship attribution	45
4.2	RCV1 studies	46
4.2.1	Studies on RCV1 on authorship attribution	47
4.3	GDELT studies	48
4.4	The guardian corpus: a case of cross-topic authorship attribution	49
4.5	Dataset selection	51
5	Our approach	53
5.1	Dataset preparation	53
5.1.1	Reuters Corpus (RCV1)	54
5.1.2	GDELT	56
5.1.3	Amazon Food Reviews (AFR)	57
5.1.4	The Guardian newspaper	57
5.2	Method's selection	59
5.2.1	Naive Approach	59
5.2.2	TeaPot	59
5.3	Features extraction	59
5.3.1	TFIDF & BOW	59
5.3.2	GridSearchCV	59
5.3.3	Doc2Vec	59
6	Results and Evaluation	55
6.1	Results single topic	55
6.1.1	RCV1 results	55
6.1.2	GDELT results	55
6.1.3	AFR results	55
6.2	Results cross topic	55
6.2.1	The Guardian results	55

7	Future works	57
8	Conclusion	59
	Bibliography	61
A	Code	71
A.1	Dataset estraction	71
A.1.1	RCV1	71
A.1.2	GDELT	71
A.2	Model	71
A.2.1	Feature extraction	71
A.2.2	Train model	71
A.2.3	Evaluation	71

OUR APPROACH

In authorship attribution problems, there is a set of candidate authors and a set of text samples in the training set covering some of the works of the authors. In the test dataset, there are sample of texts and each of them needs to be attributed to a candidate author. In the next sections, we are going to describe the experiment we carried out taking care of the chronological path of the events. Our main focus has always been on closed set authorship attribution, training with instance-based approach (i.e. extracting features by not considering the other available text samples in the training). The three milestones can be summarized as follows:

1. Dataset selection and preparation
2. Classifier method's selection
3. Features extraction

5.1 Dataset preparation

In section 4.6 we have already shown the datasets we selected. In particular, in this section we are going to show the procedures done to prepare the datasets for the next steps. For the single topic authorship attribution task we decided to select the RCV1 dataset, the dataset of 45 Victorian era book authors from the GDELT project and the dataset of amazon food reviews collected in the first decade of the 2000s. Regarding the cross domain authorship attribution task, we selected the dataset extracted from The Guardian newspaper.

5.1.1 Reuters Corpus (RCV1)

It consists of a collection of newswire stories written in English that cover four main topics: corporate/industrial (CCAT), economics (ECAT), government/social (GCAT) and markets (MCAT). We sent a request to obtain the dataset on this webpage <https://trec.nist.gov/data/reuters/reuters.html>. After few days, we gathered the RCV1 Corpus as it contains 810,000 Reuters, English Language News stories (about 2.5 GB). First of all we had to convert the dataset, that contained folders of xml files, into a big csv with author's labels and document text. Code Listing 5.1 shows the process of documents and authors extraction, using 'xml' python library. We decided to take into account this properties of the document: *text*, *title*, *headline*, *byline*, *dateline*, *lang*, *corpus_path*, *corpus_subdirectory*, *corpus_filename*.

Code Listing 5.1: Extract and Parse RCV1 XML document into csv

```
import os
import xml.etree.ElementTree as ET

for f in files:
    try:
        data_path = os.sep.join([dir_path, f])
        raw_data = open(data_path).read()
        try:
            xml_parse = ET.fromstring(raw_data)
        except:
            print(D,"/",f,"failed to parse XML.")
            continue

    def get_text(tag):
        stuff = xml_parse.find(tag)
        if stuff:
            return stuff.text
        else:
            return None

    text = "\n\n".join([str(p.text) for p in xml_parse.findall("./p")])

    title = get_text("title")
    headline = get_text("headline")
    byline = get_text("byline")
    dateline = get_text("dateline")

    #this bit got funky in the XML parse
    lang_key = [k for k in xml_parse.attrib if "lang" in k][0]
    lang = xml_parse.attrib[lang_key]
```

```

code_classes = [c.attrib["class"]
for c in xml_parse.findall("./codes")]
codes = {cc: [c.attrib["code"] for c in
xml_parse.findall("./codes[@class='%s']/code"%cc)]
for cc in code_classes}
dcs = {d.attrib["element"]: d.attrib["value"]
for d in xml_parse.findall("./dc")}

#assemble output
output = {"text": text,
"title": title,
"headline": headline,
"byline": byline,
"dateline": dateline,
"lang": lang,
"corpus_path": corpus_path,
"corpus_subdirectory": D,
"corpus_filename": f,
}

# merge and flatten the other big hashmaps
output.update(codes.items())
output.update(dcs.items())

result.append(output)
except Exception as e:
print(e)

```

The dataset was then filtered only with the documents with a *"byline"* property defined. We end up with 109'433 documents written by 2400 distinct authors. At this point, we labeled this portion of the RCV1 original dataset as the *"Full RCV1 dataset"*. In order to test and compare our approach, reproducing the testing scenario described in the previous research [59], the 10 most prolific authors were chosen from the CCAT category, and then, 50 examples per author for training and 50 examples for testing were selected randomly with no overlapping between training and testing sets. We will reference to this portion of the RCV1 dataset as the *"RCV1_10"*. In previous work [22], the authors proposed another adaptation of the RCV1 corpus for the authorship attribution task. They choose the 50 most prolific authors from the Reuters Corpus, keeping 50 examples per author for training and 50 examples per author for testing with no overlapping between them. We will refer to this corpus as the *RCV1_50*. The RCV1_10 and RCV1_50 datasets are both balanced over different authors and have their genre fixed to news. The majority of our work has been conducted on the RCV1_50,

although to compare results with previous works we will show also the same techniques applied to the RCV1_10 corpus. Table 5.1 shows the main metrics to describe these different portions of the original dataset.

Table 5.1: Main metrics to describe different portion of the dataset

Name	N# documents	N# authors	Avg docs length	Avg n# docs/author
Full RCV1 dataset	109433	2400	3061.95	45.60
RCV1_10	1000	10	3093.82	100
RCV1_50	5000	50	3251.16	100

5.1.2 GDELT

The GDELT Project is one of the largest publicly available digitized book database which has more than 3.5 million books published from 1800-2015. To decrease the bias and create a reliable dataset the following criteria have been chosen to filter out authors: English language writing authors, authors that have enough books available (at least 5), 19th century authors. With these criteria 50 authors have been selected and their books were queried through Big Query Gdelt database. The next task has been cleaning the dataset due to OCR reading problems in the original raw form. To achieve that, firstly all books have been scanned through to get the overall number of unique words and each words frequencies. While scanning the texts, the first 500 words and the last 500 words have been removed to take out specific features such as the name of the author, the name of the book and other word specific features that could make the classification task easier. After this step, we have chosen top 10, 000 words that occurred in the whole 50 authors text data corpus. The words that are not in top 10, 000 words were removed while keeping the rest of the sentence structure intact 2 . Afterwards, the words are represented with numbers from 1 to 10, 000 reverse ordered according to their frequencies. The entire book is split into text fragments with 1000 words each. We separately maintained author and book identification number for each one of them in different arrays. Text segments with less than 1000 words were filled with zeros to keep them in the dataset as well. 1000 words make approximately 2 pages of writing, which is long enough to extract a variety of features from the document. The reason why we have represented top 10, 000 words with numbers is to keep the anonymity of texts and allow researchers to run feature extraction techniques faster. Dealing with large amounts of text data can be more challenging than numerical data for some feature extraction techniques. When gathering the dataset, we decided to discard 5 authors for which their writings were not enough consistent for the authorship attribution task. We ended up with a full dataset with 53'678 documents instances, each one containing 1000 words. In order to make training's methods reliable

across dataset, we decided to select 100 documents of each authors, with a 50/50 split (i.e. 50 documents in the training set, 50 documents in the testing set, no overlapping among them). In the following sections, we will refer to this as the "*GDELT_45*". Table 5.2 shows the metrics that describe best this dataset.

Table 5.2: Main metrics to describe different portion of the GDELT dataset

Name	N# documents	N# authors	Avg docs length	Avg n# docs/author
Full GDELT dataset	53678	45	4950.61	1192.84
GDELT_45	4500	45	4911.91	100

5.1.3 Amazon Food Reviews (AFR)

This dataset consists of reviews of fine foods from amazon. The data span a period of more than 10 years, including all 500,000 reviews up to October 2012. Reviews include product and user information, ratings, and a plain text review. It also includes reviews from all other Amazon categories. We decided to consider this dataset for our experiment, because we were missing a more "everyday" example of dataset to work with. As Table 5.3 shows, among the main metrics, that the average documents length is dramatically lower than the other two dataset presented previously, providing a good challenge for us to show consistency of our method across all these different scenarios. Moreover, in order to make training's methods reliable across dataset, we decided to select 100 reviews of each customers, with a 50/50 split (i.e. 50 reviews in the training set, 50 reviews in the testing set, no overlapping among them). In the following sections, we will refer to this as the "*AFR_50*".

Table 5.3: Main metrics to describe different portion of the AFR dataset

Name	N# documents	N# authors	Avg docs length	Avg n# docs/author
Full AFR dataset	568454	256059	380.70	2.2
AFR_50	5000	50	990.45	100

5.1.4 The Guardian newspaper

Although the majority of our time and effort was focused on the first 3 single domain closed set authorship attribution task, we wanted to test our approach with a cross domain dataset. *The Guardian corpus* is composed of texts published in The Guardian daily newspaper. The majority of the corpus comprises opinion articles (comments). The newspaper describes the opinion articles using a set of tags indicating its subject. There

are eight top-level tags (World, U.S., U.K., Belief, Culture, Life&Style, Politics, Society), each one of them having multiple subtags. In order to test and compare our approach, we reproduce the testing scenario described in the previous research [60] using the Guardian corpus. The experimental scenario is as follows:

1. Select at most ten samples per author in each topic category (in Figure 5.1 the distribution of the samples per author for the Politics category after considering the restriction of ten samples per author is shown)
2. Use the samples in the Politics category as training set and train the classifier
3. Finally, test the classifier using another topic category different from Politics (four possible pairings)

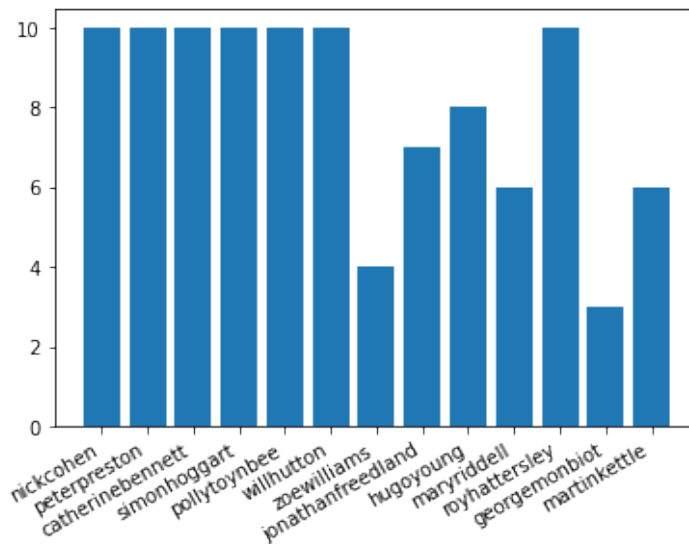


Figure 5.1: The Guardian samples distribution for the Politics topic

5.2 Method's selection

5.2.1 Naive Approach

5.2.2 TeaPot

5.3 Features extraction

5.3.1 TFIDF & BOW

5.3.2 GridSearchCV

5.3.3 Doc2Vec

BIBLIOGRAPHY

- [1] Ahmed Abbasi and Hsinchun Chen. Applying authorship analysis to extremist-group web forum messages. *IEEE Intelligent Systems*, 20(5):67–75, 2005.
- [2] Shlomo Argamon and Shlomo Levitan. Measuring the usefulness of function words for authorship attribution. In *Proceedings of the 2005 ACH/ALLC Conference*, pages 4–7, 2005.
- [3] Harald Baayen, Hans van Halteren, Anneke Neijt, and Fiona Tweedie. An experiment in authorship attribution. In *6th JADT*, volume 1, pages 69–75. Citeseer, 2002.
- [4] Sarkhan Badirli, Mary Borgo Ton, Abdulmecit Gungor, and Murat Dundar. Open set authorship attribution toward demystifying victorian periodicals. *arXiv preprint arXiv:1912.08259*, 2019.
- [5] John Burrows. ‘delta’: a measure of stylistic difference and a guide to likely authorship. *Literary and linguistic computing*, 17(3):267–287, 2002.
- [6] JK Chambers, P Trudgill, and Natalie Schilling-Estes. The handbook of language variation and change (2nd). *Victoria: Blackwell Publishing*, 2004.
- [7] Carole E Chaski. Who’s at the keyboard? authorship attribution in digital evidence investigations. *International journal of digital evidence*, 4(1):1–13, 2005.
- [8] Na Cheng, Rajarathnam Chandramouli, and KP Subbalakshmi. Author gender identification from text. *Digital Investigation*, 8(1):78–88, 2011.
- [9] Cindy Chung and James W Pennebaker. The psychological functions of function words. *Social communication*, 1:343–359, 2007.
- [10] Joachim Diederich, Jörg Kindermann, Edda Leopold, and Gerhard Paass. Authorship attribution with support vector machines. *Applied intelligence*, 19(1):109–123, 2003.
- [11] Harris Drucker, Donghui Wu, and Vladimir N Vapnik. Support vector machines for spam categorization. *IEEE Transactions on Neural networks*, 10(5):1048–1054, 1999.

- [12] Susan Dumais, John Platt, David Heckerman, and Mehran Sahami. Inductive learning algorithms and representations for text categorization. In *Proceedings of the seventh international conference on Information and knowledge management*, pages 148–155, 1998.
- [13] Neal Fox, Omran Ehmoda, and Eugene Charniak. Statistical stylometrics and the marlowe-shakespeare authorship debate. *Proceedings of the Georgetown University Roundtable on Language and Linguistics (GURT), Washington, DC, USA*, 2012.
- [14] Georgia Frantzeskou, Efstathios Stamatatos, Stefanos Gritzalis, and Sokratis Katsikas. Effective identification of source code authors using byte-level information. In *Proceedings of the 28th international conference on Software engineering*, pages 893–896, 2006.
- [15] Helena Gómez-Adorno, Juan-Pablo Posadas-Durán, Grigori Sidorov, and David Pinto. Document embeddings learned on various types of n-grams for cross-topic authorship attribution. *Computing*, 100(7):741–756, 2018.
- [16] Jack Grieve. Quantitative authorship attribution: An evaluation of techniques. *Literary and linguistic computing*, 22(3):251–270, 2007.
- [17] Abdulmecit Gungor. *Benchmarking authorship attribution techniques using over a thousand books by fifty Victorian era novelists*. PhD thesis, 2018.
- [18] H van Halteren. Linguistic profiling for authorship recognition and verification. 2004.
- [19] Graeme Hirst and Ol’ga Feiguina. Bigrams of syntactic labels for authorship discrimination of short texts. *Literary and Linguistic Computing*, 22(4):405–417, 2007.
- [20] David I Holmes. The evolution of stylometry in humanities scholarship. *Literary and linguistic computing*, 13(3):111–117, 1998.
- [21] David I Holmes and Fiona J Tweedie. Forensic stylometry: A review of the cusum controversy. *Revue Informatique et Statistique dans les Sciences Humaines*, 31(1):19–47, 1995.
- [22] John Houvardas and Efstathios Stamatatos. N-gram feature selection for authorship identification. In *International conference on artificial intelligence: Methodology, systems, and applications*, pages 77–86. Springer, 2006.
- [23] Thorsten Joachims. Making large-scale svm learning practical. Technical report, Technical report, 1998.

- [24] Thorsten Joachims et al. Transductive inference for text classification using support vector machines. In *Icml*, volume 99, pages 200–209, 1999.
- [25] Patrick Juola. *Authorship attribution*, volume 3. Now Publishers Inc, 2008.
- [26] Brett Kessler, Geoffrey Nunberg, and Hinrich Schütze. Automatic detection of text genre. *arXiv preprint cmp-lg/9707002*, 1997.
- [27] Mike Kestemont, Efstathios Stamatatos, Enrique Manjavacas, Walter Daelemans, Martin Potthast, and Benno Stein. Overview of the cross-domain authorship attribution task at pan 2019. In *CLEF (Working Notes)*, 2019.
- [28] Dmitry V Khmelev and William J Teahan. A repetition based measure for verification of text collections and for text categorization. In *Proceedings of the 26th annual international ACM SIGIR conference on Research and development in information retrieval*, pages 104–110, 2003.
- [29] Moshe Koppel and Jonathan Schler. Exploiting stylistic idiosyncrasies for authorship attribution. In *Proceedings of IJCAI’03 Workshop on Computational Approaches to Style Analysis and Synthesis*, volume 69, pages 72–80, 2003.
- [30] Moshe Koppel, Jonathan Schler, and Kfir Zigdon. Determining an author’s native language by mining a text for errors. In *Proceedings of the eleventh ACM SIGKDD international conference on Knowledge discovery in data mining*, pages 624–628, 2005.
- [31] Moshe Koppel, Jonathan Schler, Shlomo Argamon, and Eran Messeri. Authorship attribution with thousands of candidate authors. In *Proceedings of the 29th annual international ACM SIGIR conference on Research and development in information retrieval*, pages 659–660, 2006.
- [32] Moshe Koppel, Jonathan Schler, and Shlomo Argamon. Computational methods in authorship attribution. *Journal of the American Society for information Science and Technology*, 60(1):9–26, 2009.
- [33] Moshe Koppel, Jonathan Schler, and Shlomo Argamon. Authorship attribution in the wild. *Language Resources and Evaluation*, 45(1):83–94, 2011.
- [34] Moshe Koppel, Jonathan Schler, Shlomo Argamon, and Yaron Winter. The “fundamental problem” of authorship attribution. *English Studies*, 93(3):284–291, 2012.
- [35] Robert Layton, Paul Watters, and Richard Dazeley. Authorship attribution for twitter in 140 characters or less. In *2010 Second Cybercrime and Trustworthy Computing Workshop*, pages 1–8. IEEE, 2010.

- [36] Quoc Le and Tomas Mikolov. Distributed representations of sentences and documents. In *International conference on machine learning*, pages 1188–1196. PMLR, 2014.
- [37] David D Lewis and Marc Ringuette. A comparison of two learning algorithms for text categorization. In *Third annual symposium on document analysis and information retrieval*, volume 33, pages 81–93, 1994.
- [38] David Madigan, Alexander Genkin, David D Lewis, Shlomo Argamon, Dmitriy Fradkin, and Li Ye. Author identification on the large scale. In *Proceedings of the 2005 Meeting of the Classification Society of North America (CSNA)*, 2005.
- [39] Ilia Markov, Efstathios Stamatatos, and Grigori Sidorov. Improving cross-topic authorship attribution: The role of pre-processing. In *International Conference on Computational Linguistics and Intelligent Text Processing*, pages 289–302. Springer, 2017.
- [40] Yuval Marton, Ning Wu, and Lisa Hellerstein. On compression-based text classification. In *European Conference on Information Retrieval*, pages 300–314. Springer, 2005.
- [41] S Michaelson and A Morton. The qsum plot. Technical report, Internal Report CSR-3, 1990.
- [42] Tomáš Mikolov, Martin Karafiát, Lukáš Burget, Jan Černocký, and Sanjeev Khudanpur. Recurrent neural network based language model. In *Eleventh annual conference of the international speech communication association*, 2010.
- [43] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*, 2013.
- [44] Tomáš Mikolov, Wen-tau Yih, and Geoffrey Zweig. Linguistic regularities in continuous space word representations. In *Proceedings of the 2013 conference of the north american chapter of the association for computational linguistics: Human language technologies*, pages 746–751, 2013.
- [45] Leonid A Mironovsky, Alexander V Nikitin, Nina N Reshetnikova, and Nikolay V Soloviev. Graphological analysis and identification of handwritten texts. In *Computer Vision in Control Systems-4*, pages 11–40. Springer, 2018.
- [46] Tom M Mitchell. Artificial neural networks. *Machine learning*, 45:81–127, 1997.
- [47] Frederick Mosteller and David L Wallace. *Inference and disputed authorship: The Federalist*. Stanford Univ Center for the Study, 2007.

- [48] Rebekah Overdorf and Rachel Greenstadt. Blogs, twitter feeds, and reddit comments: Cross-domain authorship attribution. *Proceedings on Privacy Enhancing Technologies*, 2016(3):155–171, 2016.
- [49] Spyridon Plakias and Efstathios Stamatatos. Tensor space models for authorship identification. In *Hellenic Conference on Artificial Intelligence*, pages 239–249. Springer, 2008.
- [50] Juan-Pablo Posadas-Durán, Helena Gómez-Adorno, Grigori Sidorov, Ildar Batyrshin, David Pinto, and Liliana Chanona-Hernández. Application of the distributed document representation in the authorship attribution task for small corpora. *Soft Computing*, 21(3):627–639, 2017.
- [51] Martin Potthast, Sarah Braun, Tolga Buz, Fabian Duffhauss, Florian Friedrich, Jörg Marvin Güllow, Jakob Köhler, Winfried Löttsch, Fabian Müller, Maike Elisa Müller, et al. Who wrote the web? revisiting influential author identification research applicable to information retrieval. In *European Conference on Information Retrieval*, pages 393–407. Springer, 2016.
- [52] Joseph Rudman. The state of authorship attribution studies: Some problems and solutions. *Computers and the Humanities*, 31(4):351–365, 1997.
- [53] Conrad Sanderson and Simon Guenter. Short text authorship attribution via sequence kernels, markov chains and author unmasking: An investigation. In *Proceedings of the 2006 conference on empirical methods in natural language processing*, pages 482–491, 2006.
- [54] Upendra Sapkota, Thamar Solorio, Manuel Montes, Steven Bethard, and Paolo Rosso. Cross-topic authorship attribution: Will out-of-topic data help? In *Proceedings of COLING 2014, the 25th International Conference on Computational Linguistics: Technical Papers*, pages 1228–1237, 2014.
- [55] Upendra Sapkota, Steven Bethard, Manuel Montes, and Thamar Solorio. Not all character n-grams are created equal: A study in authorship attribution. In *Proceedings of the 2015 conference of the North American chapter of the association for computational linguistics: Human language technologies*, pages 93–102, 2015.
- [56] Roy Schwartz, Oren Tsur, Ari Rappoport, and Moshe Koppel. Authorship attribution of micro-messages. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 1880–1891, 2013.
- [57] Fabrizio Sebastiani. Machine learning in automated text categorization. *ACM computing surveys (CSUR)*, 34(1):1–47, 2002.

- [58] Efstathios Stamatatos. Author identification: Using text sampling to handle the class imbalance problem. *Information Processing & Management*, 44(2):790–799, 2008.
- [59] Efstathios Stamatatos. A survey of modern authorship attribution methods. *Journal of the American Society for information Science and Technology*, 60(3):538–556, 2009.
- [60] Efstathios Stamatatos. On the robustness of authorship attribution based on character n-gram features. *Journal of Law and Policy*, 21(2):421–439, 2013.
- [61] Sean Stanko, Devin Lu, and Irving Hsu. Whose book is it anyway? using machine learning to identify the author of unknown texts. *Machine Learning Final Projects*, 2013.
- [62] Andrew Tausz. Predicting the date of authorship of historical texts. *CS224N project*, 2011.
- [63] Antônio Theóphilo, Luís AM Pereira, and Anderson Rocha. A needle in a haystack? harnessing onomatopoeia and user-specific stylometrics for authorship attribution of micro-messages. In *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 2692–2696. IEEE, 2019.
- [64] Chris van der Lee and Antal van den Bosch. Exploring lexical and syntactic features for language variety identification. In *Proceedings of the Fourth Workshop on NLP for Similar Languages, Varieties and Dialects (VarDial)*, pages 190–199, 2017.
- [65] Vladimir N Vapnik. An overview of statistical learning theory. *IEEE transactions on neural networks*, 10(5):988–999, 1999.
- [66] Carrington B Williams. Mendenhall’s studies of word-length distribution in the works of shakespeare and bacon. *Biometrika*, 62(1):207–212, 1975.
- [67] Lili Yang, Chunping Li, Qiang Ding, and Li Li. Combining lexical and semantic features for short text classification. *Procedia Computer Science*, 22:78–86, 2013.
- [68] G Udny Yule. A test of tippett’s random sampling numbers. *Journal of the Royal Statistical Society*, 101(1):167–172, 1938.
- [69] Ying Zhao and Justin Zobel. Effective and scalable authorship attribution using function words. In *Asia Information Retrieval Symposium*, pages 174–189. Springer, 2005.

- [70] Rong Zheng, Jiexun Li, Hsinchun Chen, and Zan Huang. A framework for authorship identification of online messages: Writing-style features and classification techniques. *Journal of the American society for information science and technology*, 57(3): 378–393, 2006.
- [71] George Kingsley Zipf. Selected studies of the principle of relative frequency in language. 1932.
- [72] Sven Meyer Zu Eissen, Benno Stein, and Marion Kulig. Plagiarism detection without reference collections. In *Advances in data analysis*, pages 359–366. Springer, 2007.

