

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	13
1.1	Motivation and Problem Statement	13
1.2	Thesis Structure	13
2	Authorship attribution	15
2.1	State of the Art	16
2.2	Training's methods	17
2.2.1	Profile-based approach	18
2.2.2	Instance-based approach	19
2.3	Identify the task	20
2.3.1	Single domain vs Cross domain	20
2.3.2	Closed set vs Open set	21
2.4	Research question	21
2.4.1	Profiling the author	22
2.4.2	Finding a needle in a haystack	23
2.4.3	Verification of an author	23
3	Text characteristics analysis	25
3.1	Character Features	26
3.1.1	Affix n-grams	26
3.1.2	Word n-grams	27
3.1.3	Punctuation n-grams	27
3.2	Lexical Features	28
3.2.1	Bag of Words	29
3.2.2	Word N-grams	32
3.2.3	Vocabulary Richness	34
3.2.4	Stylometric features	35
3.2.5	Function Words	36
3.2.6	Term Frequency - Inverse Document Frequency	37
3.3	Syntactic Features	37

3.4	Semantic Features	38
3.4.1	Positivity and Negativity index	39
3.5	Application Specific Features	40
3.5.1	Vector embeddings of words (Word2Vec)	40
3.5.1.1	Skip-gram	41
3.5.1.2	Continuous Bag Of Words (CBOW)	42
3.5.2	Vector embeddings of documents (Doc2Vec)	42
4	Techniques and data collection for Authorship Attribution	45
4.1	Support Vector Machine	46
4.1.1	Support Vector Machine for authorship attribution	47
4.2	Studies on Reuters Corpus	48
4.2.1	Studies on Reuters Corpus on authorship attribution	49
4.3	GDELT studies	50
4.4	The guardian corpus: a case of cross-topic authorship attribution	51
5	Our approach	55
5.1	Dataset preparation	55
5.1.1	Reuters Corpus	56
5.1.2	GDELT	58
5.1.3	Amazon Food Reviews	59
5.1.4	The Guardian newspaper	59
5.2	Features extraction	60
5.2.1	Term Frequency - Inverse Document Frequency & Bag Of Words	61
5.2.2	Grid Search Cross Validation	63
5.2.3	Vector embeddings of documents	65
5.3	Method selection	67
5.3.1	Manual approach	67
5.3.2	Tree-based Pipeline Optimization Tool	69
6	Results and Evaluation	73
6.1	Metrics used	73
6.2	Results for single topic authorship attribution	73
6.2.1	Reuters Corpus results	73
6.2.2	GDELT Corpus results	74
6.2.3	Amazon Food Reviews Corpus results	74
6.3	The Guardian Corpus results	74
7	Future works	75

8	Conclusion	77
	Bibliography	79
A	Code	89
A.1	Dataset estraction	89
A.1.1	RCV1	89
A.1.2	GDELT	89
A.2	Model	89
A.2.1	Feature extraction	89
A.2.2	Train model	89
A.2.3	Evaluation	89

RESULTS AND EVALUATION

In this chapter we are going to show the best results obtained by setting up the training phase as described in the previous chapters. Most of the parameter tuning was done on just one dataset¹, and then the same setup was used to produce results for all our datasets.

6.1 Metrics used

6.2 Results for single topic authorship attribution

6.2.1 Reuters Corpus results

Table 6.1: Accuracy score and F1 macro score for Reuters Corpus 10 authors CCAT category.

Model	Accuracy	F1 macro
LinearSVC (combinedDFs), tfidf, simple tokenizer	0.878	0.876
LinearSVC (combinedDFs), tfidf, only-remove-quotes- tokenizer	0.908	0.906
LinearSVC (combinedDFs), tfidf, only-remove-quotes- tokenizer (threshold 1), ngram=(1,2)	0.922	0.921

Table 6.2: Accuracy score and F1 macro score for Reuters Corpus 50 authors CCAT category.

Model	Accuracy	F1 macro
LinearSVC (combinedDFs), tfidf, stock tokenizer	0.7644	0.76
LinearSVC (combinedDFs), tfidf, only-remove-quotes- tokenizer	0.7884	0.7842
LinearSVC (combinedDFs), tfidf, only-remove-quotes- tokenizer (threshold 1), ngram=(1,2)	0.7984	0.7949

Table 6.3: Accuracy score for Reuters Corpus 10 and 50 authors in the CCAT category.

Model	RCV1-10	RCV1-50
D2V words	0.8280	0.7524
Local histograms	0.8640	-
Tensor space models	0.8080	-
Character and word n-grams	0.7940	-
N-gram feature selection	-	0.7404
N-gram feature selection	-	0.7404
Our approach	0.9220	0.7984

Table 6.4: Accuracy score and F1 macro score for GDELT 45 authors.

Model	Accuracy	F1 macro
LinearSVC (combinedDFs), tfidf, stock tokenizer	0.7355	0.7090
LinearSVC (combinedDFs), tfidf, only-remove-quotes- tokenizer	0.7426	0.7173
LinearSVC (combinedDFs), d2v dmm	0.7716	0.7489

6.2.2 GDELT Corpus results

6.2.3 Amazon Food Reviews Corpus results

6.3 The Guardian Corpus results

¹The Reuters Corpus, RCV1

Table 6.5: Accuracy score and F1 macro score for Amazon Food Reviews 50 authors dataset.

Model	Accuracy	F1 macro
LinearSVC (combinedDFs), tfidf, simple tokenizer	0.7704	0.7674
LinearSVC (combinedDFs), tfidf, stock tokenizer	0.7836	0.7817
LinearSVC (combinedDFs), tfidf, only-remove-quotes-tokenizer, ngram=(1,2)	0.8388	0.8368

Table 6.6: Accuracy score and F1 macro score for The Guardian Corpus with LinearSVC (combinedDFs), tfidf, only-remove-quotes-tokenizer.

Training topic vs Test topic	Accuracy	F1 macro
Politics vs Books	0.7446	0.7640
Politics vs World	0.7560	0.7470
Politics vs Uk	0.7890	0.7010
Politics vs Society	0.8863	0.7430
Average	0.7940	0.7388

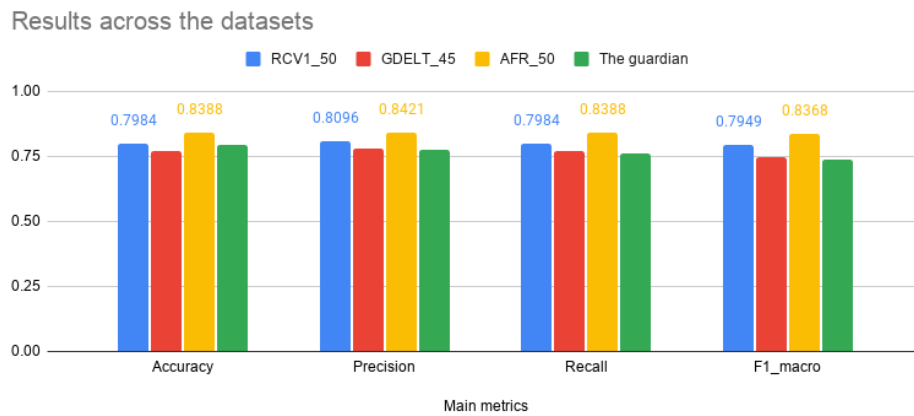


Figure 6.1: Accuracy, precision, recall and f1-macro for every dataset showing only the best result achieved for each one.

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 informaion 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] Efstathios Stamatatos. Authorship attribution using text distortion. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers*, pages 1138–1149, 2017.
- [62] 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.
- [63] Andrew Tausz. Predicting the date of authorship of historical texts. *CS224N project*, 2011.
- [64] 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.
- [65] 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.
- [66] Vladimir N Vapnik. An overview of statistical learning theory. *IEEE transactions on neural networks*, 10(5):988–999, 1999.
- [67] Carrington B Williams. Mendenhall’s studies of word-length distribution in the works of shakespeare and bacon. *Biometrika*, 62(1):207–212, 1975.
- [68] 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.
- [69] G Udny Yule. A test of tippett’s random sampling numbers. *Journal of the Royal Statistical Society*, 101(1):167–172, 1938.

- [70] Ying Zhao and Justin Zobel. Effective and scalable authorship attribution using function words. In *Asia Information Retrieval Symposium*, pages 174–189. Springer, 2005.
- [71] 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.
- [72] George Kingsley Zipf. Selected studies of the principle of relative frequency in language. 1932.
- [73] Sven Meyer Zu Eissen, Benno Stein, and Marion Kulig. Plagiarism detection without reference collections. In *Advances in data analysis*, pages 359–366. Springer, 2007.

