



Universidade Federal da Bahia
Instituto de Matemática e Estatística

Programa de Graduação em Ciência da Computação

**A DEEP LEARNING APPROACH TO
RECOGNIZE SOURCE CODE AUTHORSHIP**

Roberto Sales Caldeira

TRABALHO DE GRADUAÇÃO

Salvador
?? de dezembro de 2018

Universidade Federal da Bahia
Instituto de Matemática e Estatística

Roberto Sales Caldeira

**A DEEP LEARNING APPROACH TO RECOGNIZE SOURCE
CODE AUTHORSHIP**

*Trabalho apresentado ao Programa de Graduação em
Ciência da Computação do Instituto de Matemática e Es-
tatística da Universidade Federal da Bahia como requisito
parcial para obtenção do grau de Bacharel em Ciência da
Computação.*

Orientador: Maurício Pamplona Segundo

Salvador
?? de dezembro de 2018

Sistema de Bibliotecas - UFBA

Sales, Roberto.

A deep learning approach to recognize source code authorship / Roberto Sales Caldeira – Salvador, 2018.

15p.: il.

Orientador: Prof. Dr. Maurício Pamplona Segundo.

Monografia (Graduação) – Universidade Federal da Bahia, Instituto de Matemática e Estatística, 2018.

1. Exact algorithm. 2. Isometric embedding. 3. Maximum norm.
I. Pamplona, Maurício. II. Universidade Federal da Bahia. Instituto de Matemática e Estatística. III Título.

CDD – XXX.XX

CDU – XXX.XX.XXX

TERMO DE APROVAÇÃO

ROBERTO SALES CALDEIRA

A DEEP LEARNING APPROACH TO RECOGNIZE SOURCE CODE AUTHORSHIP

Este Trabalho de Graduação foi julgado adequado à obtenção do título de Bacharel em Ciência da Computação e aprovado em sua forma final pelo Programa de Graduação em Ciência da Computação da Universidade Federal da Bahia.

Salvador, ?? de dezembro de 2018

Prof. Dr. Maurício Pamplona Segundo
Universidade Federal da Bahia

Prof. Dr. Professor 2
Universidade Federal da Bahia

Profa. Dra. Professora 3
Universidade Federal da Bahia

DIGITE A DEDICATORIA AQUI

ACKNOWLEDGEMENTS

DIGITE OS AGRADECIMENTOS AQUI

DIGITE AQUI A CITACAO
—AUTOR (NOTA)

RESUMO

COLOQUE O RESUMO. Se preferir, crie um arquivo separado e o inclua via comando `include`.

Para evitar problemas de formato neste template (de uso geral), usamos acentuação mostrada abaixo.

`\c{c} \~{a} \'{a} \^{e} \'\{i}`

Não precisa fazer dessa forma, caso use pacotes adequados (latin1, etc.).

Palavras-chave: PALAVRAS-CHAVE.

ABSTRACT

Source code authorship identification is the task of deciding who is the author of a program given its source code. This is usually based on the analysis of previously collected samples from a set of candidate authors. There are several use cases for such a method, including attribution and detection of malicious codes, copyright infringement incident resolution, plagiarism detection, etc. As with texts in natural language, there are many distinguishing features in a piece code, like variable names, indentation style, etc. Some of these features are part of the coding style of a programmer. In this work, we investigate authorship attribution of C++ source codes based on the coding style of the authors. We also propose an end-to-end deep learning method for deciding if two source codes are from the same author, even if the involved authors are unknown to the system.

Keywords: software forensics, authorship identification, plagiarism detection

CONTENTS

Chapter 1—Introduction	1
1.1 Biometrics and Coding Style	1
1.2 Motivation	1
1.3 Applications	2
1.3.1 Plagiarism Detection	2
1.3.2 Copyright Infringement	2
1.3.3 Cyber Attack Identification	3
1.3.4 Exposing Anonymous Programmers	3
1.4 Challenges	3
1.5 Related Work	4
1.5.1 Source Code Authorship Attribution	4
1.5.2 Text Authorship Attribution	5
1.5.3 Deep Learning	5
1.6 Contribution	5
Chapter 2—Methodology	7
2.1 Problem Formulations	7
2.2 Datasets	7
2.2.1 Google Code Jam	7
2.2.2 Codeforces	8
2.3 End-to-End Deep Learning Framework	8
2.3.1 Coding Style Descriptor	8
2.3.2 Preprocessing	8
2.3.3 Models	8
2.3.3.1 Char-Level CNNs	8
Network Architecture	8
2.3.3.2 Hierarchical LSTMs	8
Network Architecture	8
2.3.4 Optimization	8
Chapter 3—Evaluation	9
Chapter 4—Conclusion	11

LIST OF FIGURES

Chapter

1

INTRODUCTION

Programmers often have to choose between tabs and spaces, between `while` loops and `for` loops, between positioning the open bracket in the next line or in the current line, etc. These are choices that can be regarded as their coding styles. Are these choices distinguishing features? In this chapter we will discuss what is style-based source code authorship identification, what challenges it poses, what has been done and what it is useful for. We will also introduce our contributions on the subject.

1.1 BIOMETRICS AND CODING STYLE

Biometrics is the field of computer vision that studies how certain human characteristics can be used to distinguish individuals. Even though physiological characteristics such as fingerprint, iris and face are probably the first ones that come to mind, the study of behavioral characteristics such as typing rhythm, voice, signature and writing style have brought new less intrusive means of distinguishing people. Even though using behavioral characteristics effectively has been made possible throughout the years, the fact that behavior is way more susceptible to change than physiological characteristics poses a big challenge.

Identifying authors of texts based on their writing styles is not a new topic (MENDENHALL, 1887). Throughout the years, computing has evolved and machine learning has reached its peek, making its way through writing style recognition. Narayanan et al. (2012) made it possible to identify the author of a text among tens of thousands of writers. It wasn't long until we were able to distinguish programmers by their coding styles (CALISKAN-ISLAM et al., 2015).

1.2 MOTIVATION

Throughout this work, we mainly consider the task of an investigator interested in deciding whether two anonymous pieces of code were authored by the same programmer or not. The actual authors of these pieces of code may be unknown to the investigator.

Also, these codes may aim to solve different problems, therefore the investigator intends to distinguish them solely based on stylistic features of such pieces.

We also consider easier scenarios where the set of possible authors are known to the investigator and labeled samples from this set are available.

We approach these problems from a deep learning perspective, training a deep neural model that can be subsequently used to resolve authorship of source codes and applying it to the different scenarios an investigator can face.

1.3 APPLICATIONS

Resolving source code authorship has a few real-world applications both in industry and academy. Although we haven't directly studied those, in this section we briefly describe a few of them. In section (...) we also describe how the experimental formulations we propose in section (...) are related to each of them.

1.3.1 Plagiarism Detection

Plagiarism can be generally defined as the unauthorized re-use of the work of another individual. Source code plagiarism is a widespread problem in academic institutions. Checking for plagiarism manually is time-consuming and not extremely effective, becoming impractical as the size of the codebase increases.

Although automatic source code plagiarism detection is a recurring and well-studied problem (MARTINS et al., 2014), the approaches consolidated by widely used tools such as MOSS (SCHLEIMER et al., 2003), Sherlock and JPlag (PRECHELT et al., 2000) are mainly based on code similarity metrics which are greatly correlated to the task the code was written to solve.

For example, let's consider the specific case of *ghostwriting*, where the code the individual claims to be of his authorship was neither written by him nor copied from a colleague, but was actually written by another person (a former student, for example). It may not be possible to compare the suspicious code with another code by the same author, since the ghostwriter may actually be unknown. On the other hand, if pieces of code of the accused party are available, it's possible to determine if his coding style matches the coding style present in the suspicious code.

Also, an analyst may strongly suspect that a piece of code a programmer claims to be of his authorship is actually not, but have no clue of who the actual author might be. This can be modeled as a binary classification problem where positive samples are other pieces of code of the same author and negative samples are pieces from unrelated programmers. We propose another approach to this problem in (...).

1.3.2 Copyright Infringement

Software forensics is the science of examining source code and binary code in order to identify, preserve, analyze and present facts and opinions about pieces of software. Although it can also be used in civil proceedings, it's most often associated with the investigation of a wide variety of computer crimes, one of which is copyright infringement.

Code correlation analysis plays an important role at copyright disputes. In this case, an analyst has labeled codes from the involved parties and the task of determining if there was infringement or not.

1.3.3 Cyber Attack Identification

Cyber attack identification is a powerful application of software forensics in cyber security. Files left behind by an intruder during a cyber attack may have just enough information for an analyst to identify who the intruder is or to relate such an attack to a previous incident. Therefore, comparing the coding style of the attacker to those of authors of previous attacks and authors of public code repositories is of great interest to the analyst.

1.3.4 Exposing Anonymous Programmers

Although there are many helpful applications of source code authorship identification, systems capable of de-anonymizing programmers pose a threat to those who want to remain anonymous, in special for anonymous open source contributors (DAUBER et al., 2017). There may be good reasons for a programmer to be anonymous, like working in a software a hostile government doesn't like.

An example of a famous open source anonymous programmer is Bitcoin's creator, Satoshi Nakamoto. For example, if we had a set of labeled codes from programmers that are likely to be Nakamoto, we could try to match their coding style to the early versions of Bitcoin, of course, assuming that Nakamoto didn't try to obfuscate his own coding style.

1.4 CHALLENGES

Although comparison metrics have proved to work well for source codes, extracting features that encode the author's style and, therefore, are independent of the task being solved have proved to be challenging. For example, features such as methods and variables names can often be misleading. This task gets even more challenging as we need to select features that are steady across different programs and capable of distinguishing between programmers. In this work, we propose an end-to-end model that solves this problem.

Also, the environment the programmer is inserted can heavily affect the difficulty of the task. For example, in projects that have a rigid style guide to be followed, much less of the programmer's own coding style might prevail. We don't study the impact of such environments in this work. Moreover, in multi-contributor projects, usually powered by VCS (version control systems), certain pieces of code can contain contributions of many authors, turning the task of relating a single author to the style present on such piece very hard. Although we believe our contributions can be applied to multi-contributor environments, we leave this for future work.

In each of the mentioned applications, the claimed author may act adversarially and try to actively modify the coding style of the program. In the ghostwriting scenario, the involved parties may act together to make the style of the code written by one to match

the other’s. During a cyber attack development, an attacker may explicitly try to hide his own coding style. In a copyright infringement, the suspect may modify the code to match his own style. In this work, adversary interference is not considered.

1.5 RELATED WORK

In this section, we describe how style-based source code authorship attribution has evolved throughout the years. We also briefly mention key works in the area of plain-text authorship attribution. While the two domains have different feature sets and different classification techniques have been applied to them, recent work in both areas are closely related. Moreover, we describe how deep learning growth might make works in these areas converge in the future.

1.5.1 Source Code Authorship Attribution

Spafford and Weeber (1993) were among the first that suggested attributing source code based on style. Even though they suggested a handful of features, they did not propose an automated method nor a thorough analysis on how those features were useful. Hayes and Offutt (2010) examined the conjecture that programmers are unique and that this uniqueness can be observed in the code they write. They conducted an experiment with programmers and graduate students, and found that programmers do have distinguishable style features which they use consistently.

Ranking approaches to source code authorship attribution were proposed by Burrows and Tahaghoghi (2007), Burrows et al. (2007, 2009) and Frantzeskou et al. (2007). Burrows and Tahaghoghi used an information retrieval technique to solve the task, obtaining token-level n -gram representations of the source codes, building an index from these representations and querying that index for programs with unknown authors. The authors of the top-ranked programs were considered the authors of the queried program. Frantzeskou et al. used byte-level n -gram features to tackle the problem. An author profile is composed of the most frequent n -grams in training data of that author. Then, the author of an unclaimed program is considered the one with the most common n -grams to this code. Both works achieved high accuracy on very small suspect sets but didn’t scale well.

Use of ASTs (abstract syntax trees) for authorship attribution was first introduced by Pellin (2000). Caliskan-Islam et al. (2015) have proposed using fuzzy ASTs and random forests to classify authorship of source code. Moreover, they proposed a coding style feature set for C/C++ source codes and a dataset for authorship attribution, based on Google Code Jam, which is a programming competition that resembles laboratory conditions. Dauber et al. (2017) showed that Caliskan-Islam et al. results could be extended to previously unexplored conditions, by adapting their techniques to work for small blocks of code of GitHub repositories.

Macdonell et al. (1999) introduced neural networks to the subject by using feed-forward neural networks and multiple discriminant analysis to attribute source codes. Bandara and Wijayarathna (2013) studied how deep neural networks could be competitive

to previous methods given enough training data.

1.5.2 Text Authorship Attribution

Plain text authorship attribution has been consistently studied over the last few decades (STAMATATOS, 2009). Most previous studies used stylometric analysis techniques for attributing authorship of literary texts (ABBASI; CHEN, 2008; NARAYANAN et al., 2012). As Internet has grown, the availability of plain text data has increased. Importance of attributing short texts, like SMS, tweets, etc. has arisen from forensics applications.

Qian et al. (2017) showed how deep learning methods can leverage from data availability, using such techniques to classify news authorship. Solorio et al. (2017) showed how CNNs (convolutional neural networks) could be used to classify authorship of short texts, like reddit comments and tweets.

1.5.3 Deep Learning

1.6 CONTRIBUTION

In this work, we introduce the concept of *coding style descriptors*, which are high dimensional representations that capture distinguishing stylistic features of a source code. We propose an end-to-end deep model that produces coding style descriptors from source code. Then, we study how the generated descriptors encode meaningful properties to the source code attribution problem by solving many of its variations.

We also introduce the Codeforces dataset for source code attribution, a C++ dataset with more than 30,000 samples extracted from Codeforces, a website specialized in holding online programming competitions. We briefly describe how the dataset was constructed and how it differs from previously published datasets.

METHODOLOGY

In this chapter, we present formulations to the many variants of the source code attribution problem in section 2.1, we describe how the Codeforces dataset was assembled in section 2.2 and present a top-down approach to how the end-to-end model was developed in section 2.3.

2.1 PROBLEM FORMULATIONS

2.2 DATASETS

The first step to develop an effective deep learning model is to gather enough training data. In this work, we decided to work with C++ source codes written in a laboratory environment – we assume the whole code is written by the author under no external style enforcement such as a style guide.

2.2.1 Google Code Jam

Although there are many public C++ laboratory datasets, the Google Code Jam¹ dataset (CALISKAN-ISLAM et al., 2015) is probably the biggest of them all. Samples from this dataset are collected from previous editions of Google Code Jam, an annual programming competition held by Google. In this competition, participants are given algorithmic tasks to be solved in a limited amount of time. As such, it’s very likely that code written by a participant manifests his own coding style.

Google Code Jam holds nearly 10 rounds every year. Most of these rounds are eliminatory. Thus, the availability of samples from less experienced participants is expected to be low. If we want to build a balanced training set not biased by the way experienced participants code, we are limited by the small amount of code less experienced participants wrote.

Although this dataset was not extensively used throughout the development phase, it was a reference for the Codeforces dataset introduced in section 2.2.2.

¹<https://codingcompetitions.withgoogle.com/codejam>

2.2.2 Codeforces

Codeforces² is a website specialized in holding online programming contests. Contest format is similar to Google Code Jam's, but they are not eliminatory. Thus, we are able to find a lot of samples from both non-experienced and experienced users.

We wrote a Python script that receives target constraints for the dataset and scrapes Codeforces for samples. Using this script, we assembled a balanced dataset with more than 30,000 samples from nearly 2,000 authors, meaning that we have around 15 samples per author. This dataset was packaged and made public³.

2.3 END-TO-END DEEP LEARNING FRAMEWORK

2.3.1 Coding Style Descriptor

The performance of machine learning methods is heavily affected by the choice of data representation. Thus, much of the effort of the machine learning community has been put into developing algorithms that transform otherwise unmanageable data into representations that can be effectively used by learning methods (BENGIO et al., 2013).

A Coding Style Descriptor (hereon referred simply as *style descriptor*) is a \mathbb{R}^n latent representation of a source code that captures its stylistic features. Ideally, style descriptors should encode everything a machine learning model needs to solve the problems posed in section 2.1. Thus, we can build simpler classifiers for these problems if we are able to build good latent representations for source codes.

Deep feed-forward networks are a natural approach to representation learning. In the remainder of this chapter, we will mainly study deep learning techniques to generate style descriptors from source codes.

2.3.2 Preprocessing

2.3.3 Models

2.3.3.1 Char-Level CNNs

Network Architecture

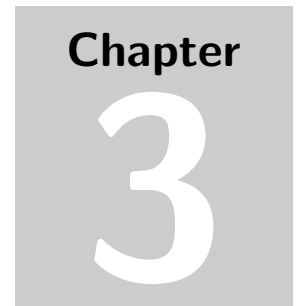
2.3.3.2 Hierarchical LSTMs

Network Architecture

2.3.4 Optimization

²<http://codeforces.com>

³link-pro-dataset



EVALUATION

Chapter

4

CONCLUSION

BIBLIOGRAPHY

- ABBASI, A.; CHEN, H. Writeprints: A stylometric approach to identity-level identification and similarity detection in cyberspace. *ACM Trans. Inf. Syst.*, ACM, New York, NY, USA, v. 26, n. 2, p. 7:1–7:29, abr. 2008. ISSN 1046-8188. Available from Internet: [⟨http://doi.acm.org/10.1145/1344411.1344413⟩](http://doi.acm.org/10.1145/1344411.1344413).
- BANDARA, U.; WIJAYARATHNA, G. Deep neural networks for source code author identification. In: *Proceedings, Part II, of the 20th International Conference on Neural Information Processing - Volume 8227*. New York, NY, USA: Springer-Verlag New York, Inc., 2013. (ICONIP 2013), p. 368–375. ISBN 978-3-642-42041-2. Available from Internet: [⟨http://dx.doi.org/10.1007/978-3-642-42042-9_46⟩](http://dx.doi.org/10.1007/978-3-642-42042-9_46).
- BENGIO, Y.; COURVILLE, A.; VINCENT, P. Representation learning: A review and new perspectives. *IEEE Trans. Pattern Anal. Mach. Intell.*, IEEE Computer Society, Washington, DC, USA, v. 35, n. 8, p. 1798–1828, ago. 2013. ISSN 0162-8828. Available from Internet: [⟨http://dx.doi.org/10.1109/TPAMI.2013.50⟩](http://dx.doi.org/10.1109/TPAMI.2013.50).
- BURROWS, S.; TAHAGHOGHI, S. M. M. Source code authorship attribution using n-grams. In: *RMIT UNIVERSITY*. [S.l.: s.n.], 2007. p. 32–39.
- BURROWS, S.; TAHAGHOGHI, S. M. M.; ZOBEL, J. Efficient plagiarism detection for large code repositories. *Softw. Pract. Exper.*, John Wiley & Sons, Inc., New York, NY, USA, v. 37, n. 2, p. 151–175, fev. 2007. ISSN 0038-0644. Available from Internet: [⟨http://dx.doi.org/10.1002/spe.v37:2⟩](http://dx.doi.org/10.1002/spe.v37:2).
- BURROWS, S.; UITDENBOGERD, A. L.; TURPIN, A. Application of information retrieval techniques for source code authorship attribution. In: *Proceedings of the 14th International Conference on Database Systems for Advanced Applications*. Berlin, Heidelberg: Springer-Verlag, 2009. (DASFAA '09), p. 699–713. ISBN 978-3-642-00886-3. Available from Internet: [⟨http://dx.doi.org/10.1007/978-3-642-00887-0_61⟩](http://dx.doi.org/10.1007/978-3-642-00887-0_61).
- CALISKAN-ISLAM, A. et al. De-anonymizing programmers via code stylometry. In: *Proceedings of the 24th USENIX Conference on Security Symposium*. Berkeley, CA, USA: USENIX Association, 2015. (SEC'15), p. 255–270. ISBN 978-1-931971-232. Available from Internet: [⟨http://dl.acm.org/citation.cfm?id=2831143.2831160⟩](http://dl.acm.org/citation.cfm?id=2831143.2831160).
- DAUBER, E. et al. Git blame who?: Stylistic authorship attribution of small, incomplete source code fragments. *CoRR*, abs/1701.05681, 2017. Available from Internet: [⟨http://arxiv.org/abs/1701.05681⟩](http://arxiv.org/abs/1701.05681).

FRANTZESKOU, G. et al. Identifying authorship by byte-level n-grams: The source code author profile (scap) method. *IJDE*, v. 6, n. 1, 2007. Available from Internet: <http://dblp.uni-trier.de/db/journals/ijde/ijde6.html\#FrantzeskouSGCH07>.

HAYES, J. H.; OFFUTT, J. Recognizing authors: an examination of the consistent programmer hypothesis. *Softw. Test., Verif. Reliab.*, v. 20, n. 4, p. 329–356, 2010. Available from Internet: <https://doi.org/10.1002/stvr.412>.

MACDONELL, S. G. et al. Software forensics for discriminating between program authors using case-based reasoning, feedforward neural networks and multiple discriminant analysis. In: *ICONIP'99. ANZIS'99 ANNES'99 ACNN'99. 6th International Conference on Neural Information Processing. Proceedings (Cat. No.99EX378)*. [S.l.: s.n.], 1999. v. 1, p. 66–71 vol.1.

MARTINS, V. T. et al. Plagiarism Detection: A Tool Survey and Comparison. In: PEREIRA, M. J. V.; LEAL, J. P.; SIMÕES, A. (Ed.). *3rd Symposium on Languages, Applications and Technologies*. Dagstuhl, Germany: Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik, 2014. (OpenAccess Series in Informatics (OASIs), v. 38), p. 143–158. ISBN 978-3-939897-68-2. ISSN 2190-6807. Available from Internet: <http://drops.dagstuhl.de/opus/volltexte/2014/4566>.

MENDENHALL, T. C. The characteristic curves of composition. *Science*, American Association for the Advancement of Science, ns-9, n. 214S, p. 237–246, 1887. ISSN 0036-8075. Available from Internet: <http://science.sciencemag.org/content/ns-9/214S/237>.

NARAYANAN, A. et al. On the feasibility of internet-scale author identification. In: *Proceedings of the 2012 IEEE Symposium on Security and Privacy*. Washington, DC, USA: IEEE Computer Society, 2012. (SP '12), p. 300–314. ISBN 978-0-7695-4681-0. Available from Internet: <https://doi.org/10.1109/SP.2012.46>.

PELLIN, B. Using classification techniques to determine source code authorship. 2000.

PRECHELT, L.; MALPOHL, G.; PHILIPPSEN, M. *JPlag: Finding plagiarisms among a set of programs*. [S.l.], 2000.

QIAN, C.; HE, T.; ZHANG, R. Deep learning based authorship identification. In: . [S.l.: s.n.], 2017.

SCHLEIMER, S.; WILKERSON, D. S.; AIKEN, A. Winnowing: Local algorithms for document fingerprinting. In: *Proceedings of the 2003 ACM SIGMOD International Conference on Management of Data*. New York, NY, USA: ACM, 2003. (SIGMOD '03), p. 76–85. ISBN 1-58113-634-X. Available from Internet: <http://doi.acm.org/10.1145/872757.872770>.

SOLORIO, T. et al. Convolutional neural networks for authorship attribution of short texts. In: *EACL*. [S.l.: s.n.], 2017.

SPAFFORD, E. H.; WEEBER, S. A. Software forensics: Can we track code to its authors? *Comput. Secur.*, Elsevier Advanced Technology Publications, Oxford, UK, UK, v. 12, n. 6, p. 585–595, out. 1993. ISSN 0167-4048. Available from Internet: [http://dx.doi.org/10.1016/0167-4048\(93\)90055-A](http://dx.doi.org/10.1016/0167-4048(93)90055-A).

STAMATATOS, E. A survey of modern authorship attribution methods. *JASIST*, v. 60, p. 538–556, 2009.