

Intrinsic plagiarism analysis

Benno Stein · Nedim Lipka · Peter Prettenhofer

Published online: 20 January 2010
© Springer Science+Business Media B.V. 2010

Abstract Research in automatic text plagiarism detection focuses on algorithms that compare suspicious documents against a collection of reference documents. Recent approaches perform well in identifying copied or modified foreign sections, but they assume a closed world where a reference collection is given. This article investigates the question whether plagiarism can be detected by a computer program if no reference can be provided, e.g., if the foreign sections stem from a book that is not available in digital form. We call this problem class *intrinsic plagiarism analysis*; it is closely related to the problem of authorship verification. Our contributions are threefold. (1) We organize the algorithmic building blocks for intrinsic plagiarism analysis and authorship verification and survey the state of the art. (2) We show how the meta learning approach of Koppel and Schler, termed “unmasking”, can be employed to post-process unreliable stylometric analysis results. (3) We operationalize and evaluate an analysis chain that combines document chunking, style model computation, one-class classification, and meta learning.

Keywords Plagiarism detection · Authorship verification · Stylometry · One-class classification

B. Stein (✉) · N. Lipka · P. Prettenhofer
Faculty of Media, Media Systems, Bauhaus-Universität Weimar, 99421 Weimar, Germany
e-mail: benno.stein@uni-weimar.de

N. Lipka
e-mail: nedim.lipka@uni-weimar.de

P. Prettenhofer
e-mail: peter.prettenhofer@uni-weimar.de

1 Problem statement

In the following, the term plagiarism refers to *text* plagiarism, i.e., the use of another author's information, language, or writing, when done without proper acknowledgment of the original source. Plagiarism detection refers to the unveiling of text plagiarism. Existing approaches to *computer-based* plagiarism detection break down this task into manageable parts:

“Given a text d and a reference collection D , does d contain a section s for which one can find a document $d_i \in D$ that contains a section s_i such that under some retrieval model \mathcal{R} the similarity $\varphi_{\mathcal{R}}$ between s and s_i is above a threshold θ ?”

Observe that research on automated plagiarism detection presumes a closed world where a reference collection D is given. Since D can be extremely large—possibly the entire indexed part of the World Wide Web—the main research focus is on efficient search technology: near-similarity search and near-duplicate detection (Brin et al. 1995; Hoad and Zobel 2003; Bernstein and Zobel 2004; Henzinger 2006; Hinton and Salakhutdinov 2006; Yang and Callan 2006), tailored indexes for near-duplicate detection (Finkel et al. 2002; Bernstein and Zobel 2004; Broder et al. 2006) or similarity hashing techniques (Kleinberg 1997; Indyk and Motwani 1998; Gionis 1999; Stein 2005, 2007). This article, however, deals with technology to identify plagiarized sections in a text if no reference collection is given. We distinguish the two analysis challenges as external and intrinsic analysis respectively. Note that human readers are able to identify plagiarism without having a reference collection at their disposal: changes between brilliant and baffling passages, or the change of person narrative give hints to multiple authorship.

1.1 Intrinsic plagiarism analysis and authorship verification

Intrinsic plagiarism analysis is closely related to authorship verification: goal of the former is to identify potential plagiarism by analyzing a document with respect to undeclared changes in writing style. Similarly, in an authorship verification problem one is given writing examples of an author A , and one is asked to determine whether or not a text with doubtful authorship is also from A . Intrinsic plagiarism analysis can be understood as a more general form of the authorship verification problem:

1. one is given a single document only, and
2. one is faced with the problem of finding the suspicious sections.

Intrinsic plagiarism analysis and authorship verification are one-class classification problems. A one-class classification problem defines a target class for which a certain number of examples exist. Objects outside the target class are called outliers, and the classification task is to tell apart outliers from target class members. Actually, the set of “outliers” can be much bigger than the target class, and an arbitrary number of outlier examples could be collected. Hence a one-class classification problem may look like a two-class discrimination problem, but there is

an important difference: members of the target class can be considered as representatives for their class, whereas one will not be able to compile a set of outliers that is representative for some kind of “non-target class”. This fact is rooted in the huge number and the diversity of possible non-target objects. Put another way, solving a one-class classification problem means to learn a concept (the concept of the target class) in the absence of discriminating features. However, in rare cases, knowledge about outliers can be used to construct representative counter examples related to the target class. Then a standard discrimination strategy can be followed.

1.2 Decision problems

Within the classical authorship verification problem the target class is comprised of writing examples of a known author A , and each piece of text written by an author B , $B \neq A$, is considered as a (style) outlier. Intrinsic plagiarism analysis is an intricate variant of authorship verification, imposing particular constraints and assumptions on the availability of writing style examples. To organize existing research we introduce the following authorship verification problems, formulated as decision problems.

1. **Problem.** AV_{EXTERNAL}

Given. A text d , written by author A , and a set of texts, $D = \{d_1, \dots, d_n\}$, written by authors \mathbf{B} , $A \notin \mathbf{B}$.

Question. Does d contain a section whose similarity to a section in d_i , $d_i \in D$, is above a threshold θ ?

2. **Problem.** AV_{FIND}

Given. A text d , allegedly written by author A .

Question. Does d contain a section written by an author B , $B \neq A$?

3. **Problem.** AV_{OUTLIER}

Given. A set of texts $D = \{d_1, \dots, d_n\}$, written by author A , and a text d , allegedly written by author A .

Question. Is d written by an author B , $B \neq A$?

The problem class AV_{EXTERNAL} corresponds to the external plagiarism analysis problem mentioned at the outset; the problem class AV_{FIND} corresponds to the general intrinsic plagiarism analysis problem, and the problem class AV_{OUTLIER} corresponds to the classical authorship verification problem. An instance π of AV_{FIND} can be reduced to m instances of AV_{OUTLIER} , $AV_{\text{FIND}} \leq_n^p AV_{\text{OUTLIER}}$, by applying a canonical chunking strategy that splits a document into m sections while asking for each section whether it forms an outlier or not. If at least one instance of AV_{OUTLIER} is answered with yes, the answer to π is yes.¹ Likewise, an instance π of AV_{OUTLIER} can be reduced to an instance of AV_{FIND} , $AV_{\text{OUTLIER}} \leq AV_{\text{FIND}}$, by

¹ The reduction \leq_n^p is in $O(|d|^2)$; within this time all possible outliers can be constructed for a document d . The reduction \leq_n^p computes the answer to AV_{FIND} from the m answers to AV_{OUTLIER} by means of a truth table π , which is a disjunction here.

simply merging d and all documents in D into a single document. The different complexity of the problem classes is reflected by the reductions $\leq_{\#}^p$ and \leq .

If the answer to an instance π of AV_{FIND} is given via a reduction of π to m AV_{OUTLIER} problems, one can try to raise the evidence of this answer by a post-processing step: from the m potential outlier sections two sets D_1 and D_2 are formed, comprising those sections that have been classified as targets into one set, and those that have been classified as outliers into the other. Again, we ask whether the documents in these two sets are written by a single author, this time applying an analysis method which takes advantage of the two sample sets, D_1 , D_2 , and which hence is more reliable than the outlier analysis. Since this decision problem is important from an algorithmic viewpoint we introduce a respective problem class:

3.' **Problem.** AV_{BATCH}

Given. A set of texts $D_1 = \{d_{1_1}, \dots, d_{1_k}\}$ written by author A , and a second set of texts, $D_2 = \{d_{2_1}, \dots, d_{2_l}\}$, allegedly written by author A .

Q. Does D_2 contain a text written by an author B , $B \neq A$?

Obviously AV_{OUTLIER} and AV_{BATCH} can be reduced to each other in polynomial time, hence $AV_{\text{OUTLIER}} \equiv AV_{\text{BATCH}}$. However, it is important to note that both reductions, $AV_{\text{FIND}} \leq_{\#}^p AV_{\text{OUTLIER}}$ and $AV_{\text{OUTLIER}} \leq AV_{\text{BATCH}}$, are constrained by a minimum text length that is necessary to perform a sensible style analysis. Experience shows that a style analysis becomes statistically unreliable for text lengths below 250 words (Stein and Meyer zu Eissen 2007).

1.3 Existing research

Authorship analysis divides into authorship *verification* problems and authorship *attribution* problems. The by far larger part of the research addresses the attribution problem: given a document d of unknown authorship and a set D of candidate authors with writing examples, and one is asked to attribute d to one author. In a verification problem (see above) one is given writing examples of an author A , and one is asked to verify whether or not a document d of unknown authorship in fact is written by A . Recent contributions to the authorship attribution problem include (Rudman 1997; Stamatatos 2001, 2007, 2009; Chaski 2005; Juola 2006; Maljutov 2006; Sanderson and Guenter 2006b); the authorship verification problem is addressed in Koppel and Schler (2004b), van Halteren (2004, 2007), Meyer zu Eissen and Stein (2006, 2007), Koppel et al. (2007), Stein and Meyer zu Eissen (2007), Stein et al. 2008 and Pavelec et al. (2008).

Several research areas are related to authorship verification, in particular: (1) stylometry, i.e., the construction of models for the quantification of writing style, text complexity, and grading level assessment, (2) outlier analysis and meta learning (Tax 2001; Tax and Duin 2001; Manevitz and Yousef 2001; Rättsch et al. 2002; Koppel and Schler 2003, 2004b, 2006), and (3) symbolic knowledge processing, i.e., knowledge representation, deduction, and heuristic inference (Russel and Norvig 1995; Stefik 1995).

In their excellent paper from 2004 Koppel and Schler give an illustrative discussion of authorship verification as a one-class classification problem (Koppel and Schler 2004b). At the same place they introduce the unmasking approach to determine whether a set of writing examples is a subset of the target class. Observe the term “set” in this connection: unmasking does not solve the one-class classification problem for a single object but requires a batch of objects all of which must stem either from the target class or not.

2 Building blocks to operationalize authorship verification

Plagiarism detection can be operationalized by decomposing a document into natural sections, such as sentences, chapters, or topically related blocks, and analyzing the variance of stylometric features for these sections. In this regard the decision problems in Sect. 1.2 are of decreasing complexity: instances of AV_{FIND} are comprised of both a selection problem (finding suspicious sections) and an AV_{OUTLIER} problem; instances of AV_{BATCH} are a restricted variant of AV_{OUTLIER} since one has the additional knowledge that all elements of a batch are (or are not) outliers at the same time.

Solving instances of AV_{FIND} involves various subtasks; Table 1 organizes them as building blocks—from left to right—following the logical text processing chain. Among others the building blocks denote alternative decomposition strategies, alternative style models, alternative classification technology, as well as post-processing options whose objective is to improve the analysis’ overall precision and recall. The table highlights those building blocks that are combined in our analysis chain; the following subsections discuss them in greater detail. Note that even with a

Table 1 Building blocks to operationalize authorship verification

Impurity assessment	Decomposition strategy	Style model construction	Outlier identification	Outlier post-processing
Document length analysis	Uniform length	Lexical character features	One-class density estimation	Heuristic voting
Genre Analysis	Structural boundaries	Lexical word features	One-class boundary estimation	Citation analysis
Analysis of issuing institution	Text element boundaries	Syntactical features	One-class reconstruction	Human inspection
	Topical boundaries	Structural features	Two-class discrimination	Unmasking
	Stylistic boundaries	Language modeling		Qsum
				Batch means

The first column lists pre-analysis methods, the second to the fourth column list the modeling and classifier methods which form the heart of a verification process, and the last column lists post-processing methods to improve the analysis quality. The highlighted building blocks indicate the employed technology of the analysis chain in this article

skillful combination and adaptation of these building blocks it is pretty difficult to end up with an analysis process comparable to the power of a human reader.

2.1 Impurity assessment

How likely is the fact that a document d contains a section of another author? We expect that the lengths, the places, and the entire fraction θ of such sections depend on particular document characteristics. Hence it makes sense to analyze the document type (paper, dissertation), its genre (novel, factual report, research, dictionary entry), but also the issuing institution (university, company, public service). Algorithmic means to reveal such information interpret document lengths, genres, and occurring named entities.

2.2 Decomposition strategy

The simplest strategy is to decompose a text d into sections s_1, \dots, s_n of uniform length; in Meyer zu Eissen and Stein (2006) the authors integrate an additional sentence detection. However, a more sensible interpretation of structural boundaries (chapters, paragraphs) is possible, which may consider special text elements like tables, formulas, footnotes, or quotations as well (Reynar 1998). Though quite difficult, the detection of topical boundaries has a significant impact on the usefulness of a decomposition (Choi 2000). In Graham et al. (2005) the authors even try to identify stylistic boundaries.

2.3 Style model construction

The statistical analysis of literary style is called stylometry, and the first ideas date back to 1851 (Holmes 1998). The automation of this task requires a quantifiable style model, and efforts in this direction became a more active research field in the 1930s (Zipf 1932; Yule 1944; Flesch 1948). In the meantime various stylometric features, also termed style markers, have been proposed. They measure writer-specific aspects like vocabulary richness (Honore 1979; Yule 1944), text complexity and understandability (Flesch 1948), or reader-specific grading levels that are necessary to understand a text (Dale and Chall 1948; Kincaid et al. 1975; Chall and Dale 1995). Note that the mentioned style features have been developed to judge longer texts, ranging from a few pages up to book size.

Style model construction must consider the decomposition strategy: different stylometric features have different strengths and also pose different constraints on text length, text genre, or topic variation. Since text plagiarism typically relates to sections that are shorter than a single page (Mansfield 2004), the decomposition of a document into sections s_1, \dots, s_n must not be too coarse, and, it is questionable which of the stylometric features will work for short sections. It should be clear that style features that employ measures like average paragraph length are not reliable in general. The authors in Meyer zu Eissen and Stein (2007) investigate the robustness of the vocabulary richness measures Yule's K , Honore's R , and the average word

frequency class. They observe that the average word frequency class can be called robust: it provides reliable results even for short sections, which can be explained with its word-based granularity. In Meyer zu Eissen and Stein (2006) connections of this type have been analyzed for the Flesch Kincaid Grade Level (1948, 1975), the Dale–Chall formula (1948, 1995), Yule’s K (1944), Honore’s R (1979), the Gunning Fog index (1952), and the averaged word frequency class (Meyer zu Eissen and Stein 2004).

Table 2 compiles an overview of important stylometric features that have been proposed so far; we distinguish between lexical features (character-based and word-based), syntactic features, and structural features. Our overview is restricted to the well-known style features and omits esoteric variants. Those features marked with an asterisk have been reported to be particularly discriminative for authorship analysis and are used within our stylometric analysis.

2.4 Outlier identification

The decomposition of a document d gives a sequence s_1, \dots, s_n of sections, for which the computation of a style model gives a sequence $\mathbf{s}_1, \dots, \mathbf{s}_n$ of feature vectors, which in turn are analyzed with respect to outliers. The identification of outliers among the s_i has to be solved on the basis of positive examples only and hence poses a one-class classification problem. Following Tax, one-class classification approaches fall into one of the following three classes (Tax 2001):

- (a) Density methods, which directly estimate the probability distributions of features for the target class. Outliers are assumed to be uniformly distributed, and, for example, Bayes’ rule can be applied to separate outliers from target class members.
- (b) Boundary methods, which avoid the estimation of the multi-dimensional density function but try to define a boundary around the set of target objects. The boundary computation is based on the distances between the objects in the target set.
- (c) Reconstruction methods come into play if prior knowledge for the generation process of target objects is available. Outliers can be distinguished from targets because of the higher reconstruction error they incur during the model fit.

The main advantage of boundary methods, namely to get by without assessing the multi-dimensional density function, can also be achieved with a density-based approach under Naive Bayes. Moreover, for our domain it is not clear how a boundary around the target set should be defined. We have also developed and analyzed reconstruction methods that rely on factor analysis and principal component analysis, but experienced difficulties due to unsatisfactory generalization behavior. Here, within our analysis chain, we resort to a one-class classifier of Type (a), which is outlined in the following.

Let S^t denote the event that a section $s \in \{s_1, \dots, s_n\}$ belongs to the target group (= not plagiarized); likewise, let S^o denote the event that s belongs to the outlier

Table 2 Compilation of important and well-known features used within a stylometric analysis. Features that are implemented within our style model are marked with an asterisk

Stylometric feature		Reference
Lexical features (character-based)	Character frequency	Zheng et al. (2006)
	Character n -gram frequency/ratio*	Kjell et al. (1994), Sanderson and Guenter (2006a), Juola (2006) and Koppel (2009)
	Frequency of special characters ('(', '&', '/', etc.)	Zheng et al. (2006)
	Compression rate	Stamatatos (2009)
Lexical features (word-based)	Average word length*	Holmes (1998) and Zheng et al. (2006)
	Average sentence length	Holmes (1998) and Zheng et al. (2006)
	Average number of syllables per word*	Holmes (1998)
	Word frequency	Mosteller and Wallace (1964), Holmes (1998) and Koppel (2009)
	Word n -grams frequency/ratio	Sanderson and Guenter (2006a)
	Number of hapax legomena	Tweedie and Baayen (1998) and Zheng et al. (2006)
	Number of hapax dislegomena	Tweedie and Baayen (1998) and Zheng et al. (2006)
	Dale–Chall index	Dale and Chall (1948) and Chall and Dale (1995)
	Flesch Kincaid grade level*	Flesch (1948) and Kincaid et al. (1975)
	Gunning Fog index*	Gunning (1952)
	Honore's R measure*	Honore (1979), Tweedie and Baayen (1998) and Zheng et al. (2006)
	Sichel's S measure	Tweedie and Baayen (1998) and Zheng et al. (2006)
	Yule's K measure*	Yule (1944), Holmes (1998), Tweedie and Baayen (1998) and Zheng et al. (2006)
	Type-token ratio	Yule (1944), Holmes (1998) and Zheng et al. (2006)
		Average word frequency class*
Syntactic features	Part-of-speech	Stamatatos (2009) and Koppel (2009)
	Part-of-speech n -gram frequency/ratio*	Koppel and Schler (2003) and Koppel (2009)
	Frequency of function words*	Mosteller and Wallace (1964), Holmes (1998), Argamon et al. (2003), Koppel and Schler (2003), Zheng et al. (2006) and Koppel (2009)
Structural features	Frequency of punctuations	Zheng et al. (2006)
	Average paragraph length	Zheng et al. (2006)
	Indentation	Zheng et al. (2006)
	Use of greetings and farewells	Zheng et al. (2006) and Stamatatos (2009)
	Use of signatures	Zheng et al. (2006) and Stamatatos (2009)

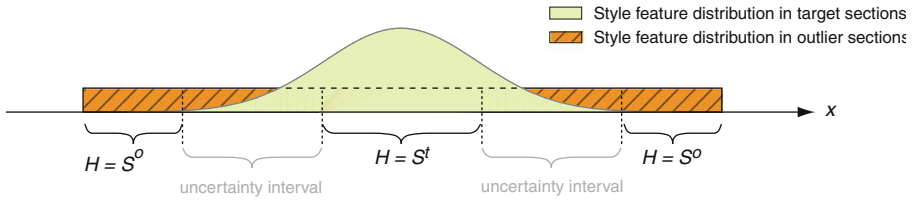


Fig. 1 Targets and outliers can be separated if they are differently distributed

group (= plagiarized). Given a document d and a single style features x , the maximum a-posteriori hypothesis $H \in \{S^t, S^o\}$ can be determined with Bayes' rule:

$$H = \operatorname{argmax}_{S \in \{S^t, S^o\}} \frac{P(x(s) | S) \cdot P(S)}{P(x(s))} \tag{1}$$

where $x(s)$ denotes the style features value for section s , and $P(x(s) | S^t)$ and $P(x(s) | S^o)$ denote the respective conditional probabilities that $x(s)$ is observed in the target group or the outlier group. Since the fraction of outliers is small compared to all sections it is sensible to estimate the $P(x(s) | S^t)$ with a Gaussian distribution; the expectation and the variance for x are estimated from $x(s_1), \dots, x(s_n)$, omitting those sections s_i that maximize or minimize $x(s_i)$. The outliers can stem from different authors, and hence the $P(x(s) | S^o)$ are estimated with a uniform distribution, following a least commitment consideration (Tax 2001). See Fig. 1 for an illustration of the assumed style feature distributions in target and outlier sections. The priors $P(S^t)$ and $P(S^o)$ correspond to $1 - \theta$ and θ respectively and require an impurity assessment (see Sect. 2.1). If no information about θ is available a uniform distribution is assumed for the priors, i.e., we resort to the maximum likelihood estimator.

Multiple style features x_1, \dots, x_m require the accounting of multiple conditional probabilities. Under the conditional independence assumption the naive Bayes approach can be applied; the accepted a-posteriori hypothesis then computes as follows:

$$H = \operatorname{argmax}_{S \in \{S^o, S^t\}} P(S) \cdot \prod_{i=1}^m P(x_i(s) | S) \tag{2}$$

For the maximum a-posteriori decision (2) only those style features x are considered whose values fall outside the uncertainty intervals (cf. Fig. 1), which are defined by 1.0 and 2.0 times the estimated standard deviation.

2.5 Outlier post-processing

The post-processing methods in Table 1 can be distinguished in knowledge-based methods and meta learning approaches. To the former count heuristic voting, citation analysis, and human inspection. Heuristic voting, which is applied here, is the estimation and use of acceptance and rejection thresholds based on the number of classified outlier sections. Meta learning is brought into play if from the solution

of several AV_{OUTLIER} problems two sets D_1 (sections labeled as targets) and D_2 (sections labeled as outliers) are formed, obtaining this way an instance of the AV_{BATCH} problem. Possible meta learning approaches are:

- (a) Unmasking (Koppel and Schler 2004b), which is a representative of what Tax terms “reconstruction method” (Tax 2001); it measures the increase of a sequence of reconstruction errors, starting with a good reconstruction which then is successively impaired.
- (b) The Qsum heuristic (Morton and Michaelson 1990; Hilton and Holmes 1993), which compares the growth rates of two cumulative sums over a sequence of sentences. Basis for the sums are the deviations from the mean sentence length and the deviations of function word frequencies.
- (c) Batch means, which is applied within the analysis of simulation data in order to detect the end of a transient phase. For a series of values the variance development of the sample mean is measured while the sample size is successively increased.

Unmasking has been successfully applied to solve instances of AV_{BATCH} (Sanderson and Guenter 2006b; Koppel and Schler 2004a; Koppel et al. 2007; Surdulescu 2004). The robustness of the approach is also reported by Kacmarcik and Gamon who develop methods for obfuscating document stylometry in order to preserve author anonymity (Kacmarcik and Gamon 2006). Since unmasking is a building block in our analysis chain it is explained in greater detail now. The use of unmasking for intrinsic plagiarism analysis was proposed in Stein and Meyer zu Eissen (2007), who consider a style outlier analysis as a heuristic to compile a potentially plagiarized and sufficiently large auxiliary document.

Recall that the set D_1 (targets) is attributed to author A , while the authorships of the sections in D_2 (outliers) is considered as unsettled. With unmasking we seek further evidence for the hypothesis whether a text in D_2 is written by an author B , $B \neq A$. At first, D_1 and D_2 are represented under a reduced vector space model, designated as \mathbf{D}_1 and \mathbf{D}_2 . As an initial feature set the 250 words with the highest relative frequency in $D_1 \cup D_2$ are chosen. Unmasking then happens in the following steps (see Fig. 2):

1. *Model Fitting*. Training of a classifier that separates \mathbf{D}_1 from \mathbf{D}_2 . In Koppel and Schler (2004b) the authors implement a tenfold cross-validation experiment with a linear kernel SVM to determine the achievable accuracy.

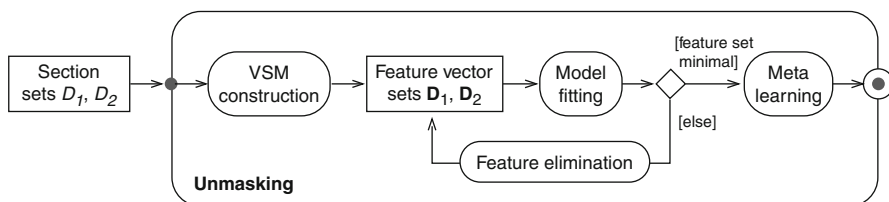


Fig. 2 Given are two sets of sections D_1 and D_2 , allegedly written by a single author. Unmasking measures the separability of \mathbf{D}_1 versus \mathbf{D}_2 when the style model is successively impaired

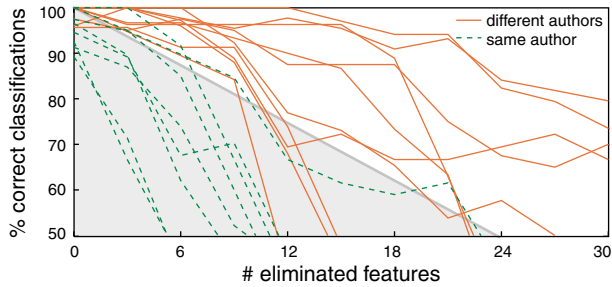


Fig. 3 Unmasking at work: each line corresponds to a comparison of two papers. A *solid red line* belongs to papers of two different authors; a *dashed green line* belongs to papers of the same author

2. *Impairing*. Elimination of the most discriminative features with respect to the model obtained in Step 1; construction of new collections \mathbf{D}_1 , \mathbf{D}_2 , which now contain impaired representations. Koppel and Schler (2004b) reports on convincing results by eliminating the six most discriminating features. This heuristic depends on the section length which in turn depends on the length of d .
3. Go to Step 1 until the feature set is sufficiently reduced. About 5–10 iterations are typical.
4. *Meta Learning*. Analyze the degradation in the quality of the model fitting process: if after the last impairing step the sets \mathbf{D}_1 and \mathbf{D}_2 can still be separated with a small error, assume that d_1 and d_2 stem from different authors. Figure 3 shows a characteristic plot where unmasking is applied to short papers of 4–8 pp.

The rationale of unmasking: Two sets of sections, D_1 , D_2 , constructed from two different documents d_1 and d_2 of the same author can be told apart easily if a vector space model (VSM) retrieval model is chosen. The VSM considers all words in $d_1 \cup d_2$, and hence it includes all kinds of open class and closed class word sets. If only the 250 most-frequent words are selected, a large fraction of them will be function words and stop words.² Among these 250 most-frequent words a small number does the major part of the discrimination job; these words capture topical differences, differences that result from genre, purpose, or the like. By eliminating them, one approaches step by step the distinctive and subconscious manifestation of an author's writing style. After several iterations the remaining features are not powerful enough to discriminate two documents of the same author. But, if d_1 and d_2 stem from two different authors, the remaining features will still quantify significant differences between \mathbf{D}_1 and \mathbf{D}_2 .

3 Analysis

This section reports on the performance of the operationalized analysis chain. Figure 4 gives an illustration: the top row shows documents with original sections

² Function words and stop words are not disjoint sets: most function words in fact are stop words; however, the converse does not hold.

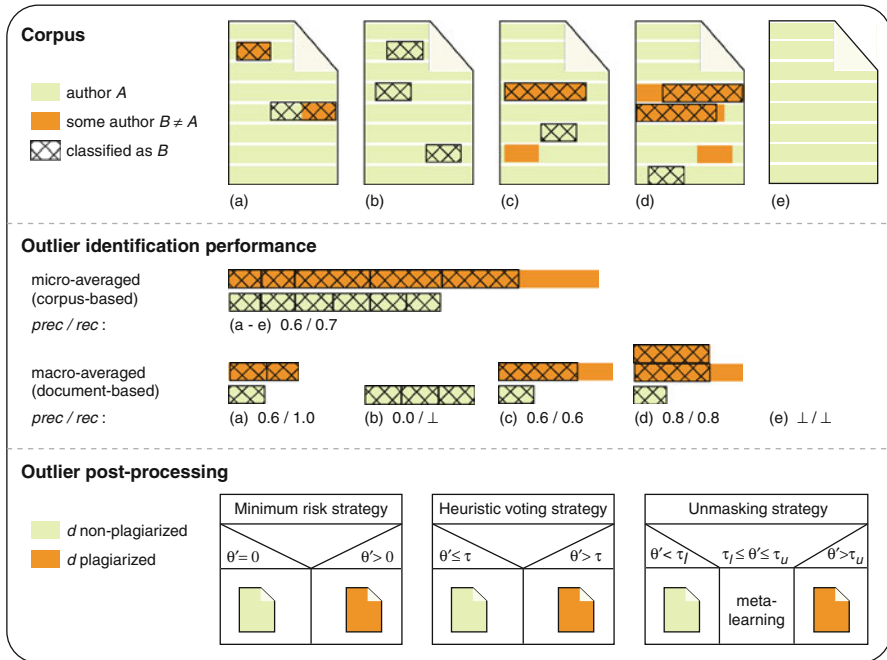


Fig. 4 Illustration of the analysis chain. *Top*: corpus with five documents of author A, containing sections of some author $B \neq A$. *Middle*: micro- and macro-averaged analysis of the outlier identification performance. *Bottom*: outlier post-processing according to three alternative strategies; θ' denotes the fraction of sections per document that are classified as outliers

(green), plagiarized sections (red), and sections spotted by the classifier (hashed); the middle row shows the micro- and macro-averaged outlier classification performance; the bottom row shows three alternative post-processing strategies. These strategies differ with respect to the interpretation of the fraction θ' of sections per document that are classified as outliers: under the minimum risk strategy a document d is considered as plagiarized if at least one outlier section is spotted, under the heuristic voting strategy θ' is compared to a threshold τ , and under the unmasking strategy meta learning is applied if θ' falls into an uncertainty interval. The remainder of this section gives particulars.

3.1 Corpus

To run analyses on a large scale one has to resort to artificially plagiarized documents. Here, we use a subset of the corpus that has been constructed for the intrinsic plagiarism analysis task of the PAN'09 competition (Potthast et al. 2009). The PAN'09 corpus comprises about 3,000 generated cases of intrinsic plagiarism—more precisely: cases of style contamination—exhibiting varying degrees of obfuscation. The corpus is based on books from the English part of the Project Gutenberg and contains mainly narrative text. Sections of varying length, ranging

Table 3 Selected summary statistics of the four test collections

Collection	No. of documents		No. of sections (total)		No. of sections (avg.)		Impurity θ (avg.)
	Plag.	Non-plag.	Plag.	Non-plag.	Plag.	Non-plag.	
1	231	231	2,067	44,316	4.5	96	0.09
2	178	178	451	9,560	1.3	27	0.09
3	178	178	4,744	21,896	13.3	62	0.30
4	188	188	1,871	7,814	5.0	21	0.33

The statistics of the columns 2–5 are per collection and consider both the plagiarized and the non-plagiarized documents; the statistics of the columns 6–8 are per document; the statistics of the columns 6–7 consider both the plagiarized and the non-plagiarized documents, whereas column 8 considers only the plagiarized documents of a collection

from a few sentences up to many pages, are inserted into other documents according to heuristic placement rules. In addition, obfuscation of the inserted sections is performed by replacing, shuffling, deleting, or adding words.³

For our experiments the documents of the PAN'09 corpus are uniformly decomposed into candidate sections of 5,000 characters; each candidate section s in turn is categorized as being either non-plagiarized, if s contains no word from an inserted section, or plagiarized, if s consists to more than 50% of an inserted section. Otherwise s is discarded and excluded from further investigations. Documents with less than seven sections are removed from the corpus because they are considered to be too short for a reliable stylometric analysis.

In order to study the effect of document length and impurity on the performance of our analysis chain, four disjoint collections are compiled. For this purpose two levels of document lengths are introduced (short versus long) and combined with two levels of impurity (light versus strong). Short documents consist of less than 250,000 characters, which corresponds to approximately 40,000 words. The impurity θ of a document is defined as the portion of plagiarized characters, i.e., characters that belong to an inserted section. A document is considered to have a light impurity if $\theta \leq 0.15$; it has a strong impurity if $\theta > 0.15$. Finally, the number of plagiarized documents per collection is set to 50%. The resulting test collections exhibit varying degrees of difficulty, both in terms of training data scarcity (document length) and class imbalance (impurity). We number the collections according to their level of difficulty and show selected summary statistics in Table 3.

3.2 Performance of outlier identification

Outlier identification is addressed with the density estimation method as described in Sect. 2.4. To capture a broad range of writing styles a diverse set of stylometric features is employed, belonging to three of the four categories introduced in Sect. 2.3: lexical character features, lexical word features, and syntactical features. Among the employed stylometric features are the classical measures for vocabulary

³ The corpus can be downloaded at <http://www.webis.de/research/corpora>.

Table 4 Stylometric features ranked by their F -measure performance in a style outlier detection task. The classification decision is given by the maximum a-posterior hypothesis from Eq. 1

Stylometric feature	F -measure
Flesch reading ease score	0.208
Average number of syllables per word	0.205
Frequency of term: of	0.192
Noun-verb-noun tri-gram	0.189
Noun-noun-verb tri-gram	0.182
Verb-noun-noun tri-gram	0.179
Gunning fog index	0.179
Yule's K measure	0.176
Flesch kincaid grade level	0.175
Average word length	0.173
Noun-preposition-propemnoun tri-gram	0.173
Honore's R measure	0.165
Average word length	0.165
Average word frequency class	0.162
Consonant-vowel-consonant tri-gram	0.154
Frequency of term: is	0.151
Noun-noun-coordinatingconjunction tri-gram	0.150
Nounplural-preposition-determiner tri-gram	0.149
Determiner-nounplural-preposition tri-gram	0.148
Consonant-vowel-vowel tri-gram	0.146
Verb-noun-verb tri-gram	0.146
Vowel-vowel-consonant tri-gram	0.146
Frequency of term: the	0.141
Determiner-noun-preposition tri-gram	0.139
Frequency of term: been	0.136
Noun-noun-noun tri-gram	0.134
Noun-preposition-determiner tri-gram	0.133
Vowel-vowel-vowel tri-gram	0.129
Noun-preposition-noun	0.128
Verb-preposition-determiner tri-gram	0.127

richness, text complexity, as well as stylometric features that have been reported to be particularly discriminative for authorship analysis, such as character n -grams and the frequency of function words (see Table 2). To capture syntactic variations in writing style, part-of-speech information in the form of part-of-speech trigrams is exploited; the tagging is done with the probabilistic part-of-speech tagger QTAG.

Table 4 shows the top 30 stylometric features with respect to their discriminative power; the F -Measure-value pertains to the outlier class and is computed as micro-averaged mean over the four collections. The decision whether or not a section is classified as an outlier is given by the maximum a-posteriori hypothesis of the univariate model in Eq. 1. Note that this ranking serves merely for illustration purposes and is not used for feature selection: the outlier analysis in the analysis

Table 5 Performance of the one-class classifier. The target class relates to sections of author A ; the outlier class relates to sections of foreign authors $B \neq A$

Collection	Target class			Outlier class		
	Prec	Rec	F	Prec	Rec	F
1	0.98	0.91	0.94	0.20	0.52	0.29
2	0.89	0.90	0.89	0.34	0.32	0.33
3	0.98	0.64	0.77	0.10	0.78	0.18
4	0.89	0.64	0.74	0.27	0.64	0.38

chain is based on the multivariate use of all stylometric features. For each document in a collection an individual style classifier according to Eq. 2 is constructed and applied to each section of that document. The correctness of each classification decision is pooled over all documents. Table 5 summarizes the achieved classification results in terms of micro-averaged F -Measure for both the outlier class and the target class.

Recall that the four collections are compiled in a way that sections with less than 50% plagiarism are discarded. If all sections with less than 90% plagiarism are discarded, the precision of the outlier class is unchanged, but its recall increases by 9% on average over all collections. On the other hand, if sections with less than 50% plagiarism are kept, the precision and the recall of the outlier class decrease by 4% on average.

3.3 Performance of meta learning

To illustrate the performance of the unmasking approach we evaluate the meta learner that is used in Step 4 of the unmasking procedure. Unmasking is parameterized as follows: documents are represented under the term frequency vector space model, defined by the 500 most frequent words of the input document sets, without applying stemming or stop wording. In each iteration i of 30 unmasking iterations the best 10 features according to the information gain heuristic are removed and the classification accuracy, acc_i , of a linear kernel SVM is computed, based on fivefold cross validation.

In practice the distribution of the outlier and target class is extremely unbalanced. In order to correct this class imbalance, the outlier class is over-sampled. Here, the SMOTE approach is used to create new, synthetic instances of the outlier class by interpolating between the original instances (Chawla et al. 2002). A meta learner is trained with vectors each of which comprising the following elements: the acc -values of iteration i , the Δ - acc -values to iteration $i - 1$, the Δ - acc -values to iteration $i - 2$, and a class label “plagiarized” or “non-plagiarized”. This meta learner is also realized as a linear kernel SVM; Table 6 reports on its performance.

The unmasking approach of Koppel and Schler decides for two sets of documents whether or not all documents stem from a single author. If both sets belong to the same author the associated unmasking curve drops away (cf. the dashed green lines in Fig. 3). This fact is exploited within our analysis chain in order to reduce the

Table 6 Evaluation of the unmasking meta learner

Collection	Non-plagiarized documents			Plagiarized documents		
	Prec	Rec	F	Prec	Rec	F
1	0.78	0.86	0.82	0.82	0.73	0.77
2	0.77	0.88	0.82	0.48	0.30	0.37
3	0.95	0.94	0.95	0.94	0.95	0.95
4	0.70	0.69	0.70	0.68	0.70	0.69

Setting: tenfold cross validated with 100 plagiarized documents and 100 non-plagiarized documents drawn randomly from the corresponding collection

number of misclassified non-plagiarized documents, which are caused by the insufficient precision of the one-class classifier.

3.4 Performance of the analysis chain

We evaluate three strategies, from naive to sophisticated, to solve AV_{FIND} for a document d . Under the minimum risk strategy d is classified as plagiarized if at least one style outlier has been announced for d . Under the heuristic voting strategy d is classified as plagiarized if the detected fraction of outlier text is above a threshold τ . Under the unmasking strategy d is classified as plagiarized if the detected fraction of outlier text is above an upper threshold τ_u ; d is classified as non-plagiarized if the detected fraction of outlier text is below a lower threshold τ_l ; for all other cases unmasking is applied. Note that the values for τ , τ_u , and τ_l are collection-dependent. In our experiments τ and τ_l are fitted to the averaged impurities of the collections, while τ_u is chosen overly optimistic. Table 7 summarizes the results: the minimum risk strategy classifies all documents as plagiarized because of the imprecision of the outlier detection, which claims at least one section in each document as outlier. Heuristic voting and unmasking consider the outlier detection characteristic. A main observation is that especially unmasking can be used to substantially increase the precision when solving instances of AV_{FIND}.

Table 7 Overall performance of the analysis chain. Performance of the solution of the AV_{FIND} problem under different strategies: minimum risk (columns 2–4), heuristic voting (columns 5–8), and unmasking (columns 9–12). Maximum precision values are shown bold

Collection	Minimum risk			Heuristic voting			Unmasking				
	Prec	Rec	F	τ	Prec	Rec	F	$(\tau_l; \tau_u)$	Prec	Rec	F
1	0.50	1.00	0.66	0.1	0.55	0.57	0.63	(0.1; 0.5)	0.83	0.50	0.62
2	0.50	1.00	0.66	0.1	0.50	1.00	0.66	(0.1; 0.5)	0.66	0.57	0.67
3	0.50	1.00	0.66	0.2	0.69	0.30	0.42	(0.2; 0.8)	0.72	0.30	0.43
4	0.50	1.00	0.66	0.2	0.52	0.97	0.68	(0.2; 0.8)	0.98	0.60	0.74

4 Summary

Intrinsic plagiarism detection is the spotting of sections with undeclared writing style changes in a text document. Intrinsic plagiarism detection is a one-class classification problem that cannot be tackled with a single technique but requires the combination of algorithmic and statistical building blocks. Our article provides an overview of these building blocks and presents ideas to operationalize analysis chains that cope with the intrinsic plagiarism challenge.

Intrinsic plagiarism detection and authorship verification are two sides of the same coin. This fact is explained in this article, and, in order to organize existing research and to work out the intricate difficulties between problem variants, we introduce four problem classes for authorship verification problems. We propose and implement an analysis chain that integrates document chunking, style model computation, style outlier identification, and outlier post-processing. Style outlier identification is unreliable, among others because it is difficult to quantify style and to spot style changes in short sections. Since we feel that plagiarism detection technology should avoid the announcement of wrongly claimed plagiarism at all costs, we propose to post-process the results of the outlier identification step. We employ the unmasking technology for this purpose, which has been developed to settle the authorship for a text in question—if sufficient sample text is at one's disposal. The combination of outlier identification with unmasking entails a significant improvement of the precision (see Table 7 for details). However, we see different places and room to improve certain building blocks in the overall picture, among others: knowledge-based chunking, better style models, multivariate one-class classification, and bootstrapping for outlier identification.

References

- Argamon, S., Šarić, M., & Stein, S. S. (2003). Style mining of electronic messages for multiple authorship discrimination: First results. In *KDD '03: Proceedings of the ninth ACM SIGKDD international conference on knowledge discovery and data mining* (pp. 475–480). New York, NY, USA: ACM. ISBN 1-58113-737-0. doi:10.1145/956750.956805.
- Bernstein, Y., & Zobel, J. (2004). A scalable system for identifying co-derivative documents. In A. Apostolico & M. Melucci (Eds.), *Proceedings of the string processing and information retrieval symposium (SPIRE)* (pp. 55–67). Padova, Italy: Springer. Published as LNCS 3246.
- Brin, S., Davis, J., & Garcia-Molina, H. (1995). Copy detection mechanisms for digital documents. In *SIGMOD '95* (pp. 398–409). New York, NY, USA: ACM Press. ISBN 0-89791-731-6.
- Broder, A. Z., Eiron, N., Fontoura, M., Herscovici, M., Lempel, R., McPherson, J., et al. (2006). Indexing shared content in information retrieval systems. In *EDBT '06* (pp. 313–330).
- Chall, J. S., & Dale, E. (1995). *Readability revisited: The new Dale–Chall readability formula*. Cambridge, MA: Brookline Books.
- Chaski, C. E. (2005). Who's at the keyboard? authorship attribution in digital evidence investigations. *IJDE*, 4(1), 1–14.
- Chawla, N. V., Bowyer, K. W., Kegelmeyer, P. W. (2002). Smote: Synthetic minority over-sampling technique. *Journal of Artificial Intelligence Research* 16, 321–357.
- Choi, F. Y. Y. (2000). Advances in domain independent linear text segmentation. In *Proceedings of the first conference on North American chapter of the association for computational linguistics* (pp. 26–33). San Francisco, CA, USA: Morgan Kaufmann Publishers Inc.
- Dale, E., & Chall, J. S. (1948). A formula for predicting readability. *Educational Research Bulletin*, 27, 11–20.

- Finkel, R. A., Zaslavsky, A., Monostori, K., & Schmidt, H. (2002). Signature extraction for overlap detection in documents. In *Proceedings of the 25th Australian conference on Computer science* (pp. 59–64). Australian Computer Society, Inc. ISBN 0-909925-82-8.
- Flesch, R. (1948). A new readability yardstick. *Journal of Applied Psychology*, 32, 221–233.
- Gionis, A., Indyk, P., & Motwani, R. (1999). Similarity search in high dimensions via hashing. In *Proceedings of the 25th VLDB conference Edinburgh, Scotland* (pp. 518–529).
- Graham, N., Hirst, G., & Marthi, B. (2005). Segmenting a document by stylistic character. *Natural Language Engineering*, 11(4), 397–415. Supersedes August 2003 workshop version.
- Gunning, R. (1952). *The technique of clear writing*. New York: McGraw-Hill.
- Henzinger, M. (2006). Finding near-duplicate web pages: A large-scale evaluation of algorithms. In *SIGIR '06: Proceedings of the 29th annual international ACM SIGIR conference on research and development in information retrieval* (pp. 284–291). New York, NY, USA: ACM Press. ISBN 1-59593-369-7. doi:10.1145/1148170.1148222.
- Hilton, M. L., & Holmes, D. I. (1993). An assessment of cumulative sum charts for authorship attribution. *Literary and Linguistic Computing*, 8(2), 73–80.
- Hinton, G. E., & Salakhutdinov, R. R. (2006). Reducing the dimensionality of data with neural networks. *Science*, 313, 504–507.
- Hoad, T. C., & Zobel, J. (2003). Methods for identifying versioned and plagiarised documents. *American Society for Information Science and Technology*, 54(3), 203–215.
- Holmes, D. I. (1998). The evolution of stylometry in humanities scholarship. *Literary and Linguistic*, 13(3), 111–117. doi:10.1093/llc/13.3.111.
- Honore, A. (1979). Some simple measures of richness of vocabulary. *Association for Literary and Linguistic Computing Bulletin*, 7(2), 172–177.
- Indyk, P., & Motwani, R. (1998). Approximate nearest neighbor—Towards removing the curse of dimensionality. In *Proceedings of the 30th symposium on theory of computing* (pp. 604–613).
- Juola, P. (2006). Authorship attribution. *Foundation Trends Information Retrieval* 1(3), 233–334, ISSN 1554-0669. doi:10.1561/1500000005.
- Kacmarcik, G., & Gamon, M. (2006). Obfuscating document stylometry to preserve author anonymity. In *Proceedings of the COLING/ACL on main conference poster sessions* (pp. 444–451). Morristown, NJ, USA: Association for Computational Linguistics.
- Kincaid, J., Fishburne, R. P., Rogers, R. L., & Chissom, B. S. (1975). Derivation of new readability formulas (automated readability index, fog count and flesch reading ease formula) for navy enlisted personnel. Research branch report 8–75. Millington TN: Naval Technical Training US Naval Air Station.
- Kjell, B., Woods Addison, W., & Frieder, O. (1994). Discrimination of authorship using visualization. *Information Processing and Management*, 30(1), 141–150. ISSN 0306-4573. doi:10.1016/0306-4573(94)90029-9.
- Kleinberg, J. (1997). Two algorithms for nearest-neighbor search in high dimensions. In *STOC '97: Proceedings of the twenty-ninth annual ACM symposium on theory of computing*.
- Koppel, M., & Schler, J. (2003). Exploiting stylistic idiosyncrasies for authorship attribution. In *Proceedings of IJCAI'03 workshop on computational approaches to style analysis and synthesis*. Mexico: Acapulco.
- Koppel, M., & Schler, J. (2004a). Authorship verification as a one-class classification problem. In *ICML '04: Proceedings of the twenty-first international conference on Machine learning* (pp. 62). New York, NY, USA: ACM. ISBN 1-58113-828-5. doi:10.1145/1015330.1015448.
- Koppel, M., & Schler, J. (2004b). Authorship verification as a one-class classification problem. In *Proceedings of the 21st international conference on machine learning*. Banff, Canada: ACM Press.
- Koppel, M., Schler, J., Argamon, S., & Messeri, E. (2006). Authorship attribution with thousands of candidate authors. In *Proceedings of the 29th annual international ACM SIGIR conference on research and development in information retrieval* (pp. 659–660). New York, NY, USA: ACM. ISBN 1-59593-369-7. doi:10.1145/1148170.1148304.
- Koppel, M., Schler, J., & Bonchek-Dokow, E. (2007). Measuring differentiability: Unmasking pseudonymous authors. *Journal of Machine Learning Research*, 8, 1261–1276. ISSN 1533-7928.
- Koppel, M., Schler, J., & Argamon, S. (2009). Computational methods in authorship attribution. *Journal of the American Society for Information Science and Technology*, 60(1), 9–26.
- Malyutov, M. B. (2006). Authorship attribution of texts: A review. *Lecture Notes in Computer Science*, 2063, 362–380.

- Manevitz, L. M., & Yousef, M. (2001). One-class SVMs for document classification. *Journal of Machine Learning Research*, 2, 139–154.
- Mansfield, J. S. (2004). Textbook plagiarism in psy101 general psychology: incidence and prevention. In *Proceedings of the 18th annual conference on undergraduate teaching of psychology: Ideas and innovations*. New York, USA: SUNY Farmingdale.
- Meyer zu Eissen, S., & Stein, B. (2004). Genre classification of web pages: User study and feasibility analysis. In S. Biundo, T. Frühwirth, & G. Palm (Eds.), *KI 2004: Advances in artificial intelligence*, vol. 3228 LNAI of *Lecture Notes in artificial intelligence* (pp. 256–269). Berlin Heidelberg New York: Springer. ISBN 0302-9743.
- Meyer zu Eissen, S., & Stein, B. (2006). Intrinsic plagiarism detection. In M. Lalmas, A. MacFarlane, S. M. Rüger, A. Tombros, T. Tsirikika, & A. Yavlinsky (Eds.), *Proceedings of the European conference on information retrieval (ECIR 2006)*, vol. 3936 of *Lecture Notes in Computer Science* (pp. 565–569). New York: Springer. ISBN 3-540-33347-9.
- Meyer zu Eissen, S., Stein, B., & Kulig, M. (2007). Plagiarism detection without reference collections. In R. Decker & H. J. Lenz (Eds.), *Advances in data analysis* (pp. 359–366). New York: Springer. ISBN 978-3-540-70980-0.
- Morton, A. Q., & Michaelson, S. (1990). The qsum plot. Technical report, University of Edinburgh.
- Mosteller, F., & Wallace, D. L. (1964). *Inference and disputed authorship: Federalist papers*. Reading, MA: Addison-Wesley Educational Publishers Inc, 1964. ISBN 0201048655.
- Pavelec, D., Oliveira, L. S., Justino, E. J. R., & Batista, L. V. (2008). Using conjunctions and adverbs for author verification. *Journal of UCS*, 14(18), 2967–2981.
- Potthast, M., Eiselt, A., Stein, B., Barrón Cedeño, A., & Rosso, P. (Eds.). (2009). *Webis at Bauhaus-Universität Weimar and NLEL at Universidad Polytécnica de Valencia*. PAN Plagiarism Corpus 2009 (PAN-PC-09). <http://www.webis.de/research/corpora>.
- Rätsch, G., Mika, S., Schölkopf, B., & Müller, K.-R. (2002). Constructing boosting algorithms from SVMs: An application to one-class classification. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 24(9), 1184–1199. ISSN 0162-8828. doi:10.1109/TPAMI.2002.1033211.
- Reynar, J. C. (1998). *Topic segmentation: Algorithms and applications*. Ph.D. thesis, University of Pennsylvania.
- Rudman, J. (1997). The state of authorship attribution studies: Some problems and solutions. *Computers and the Humanities*, 31, 351–365.
- Russel, S. J., & Norvig, P. (1995). *Artificial intelligence: A modern approach*. Englewood Cliffs, NJ: Prentice-Hall.
- Sanderson, C., & Guenter, S. (2006a). On authorship attribution via markov chains and sequence kernels. In *Pattern recognition, 2006. ICPR 2006. 18th international conference on* (vol. 3, pp. 437–440). doi:10.1109/ICPR.2006.899.
- Sanderson, C., & Guenter, S. (2006b). 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* (pp. 482–491). URL <http://acl.ldc.upenn.edu/W/W06/W06-1657.pdf>.
- Stamatatos, E. (2007). Author identification using imbalanced and limited training texts. In A. M. Tjoa & R. R. Wagner (Eds.), *18th international conference on database and expert systems applications (DEXA 07)* (pp. 237–241). IEEE, September 2007. ISBN 0-7695-2932-1. doi: 10.1109/DEXA.2007.37.
- Stamatatos, E. (2009). A survey of modern authorship attribution methods. *Journal of American Society for Information Science & Technology*, 60(3), 538–556. ISSN 1532-2882. doi:10.1002/asi.v60.3.
- Stamatatos, E., Fakotakis, N., & Kokkinakis, G. (2001). Computer-based authorship attribution without lexical measures. *Computers and the Humanities*, 35, 193–214.
- Stefik, M. (1995). *Introduction to knowledge systems*. San Mateo, CA, USA: Morgan Kaufmann.
- Stein, B. (2005). Fuzzy-fingerprints for text-based information retrieval. In K. Tochtermann & H. Maurer (Eds.), *Proceedings of the 5th international conference on knowledge management (I-KNOW 05)*, Graz, *Journal of Universal Computer Science* (pp. 572–579). Know-Center.
- Stein, B. (2007). Principles of hash-based text retrieval. In C. Clarke, N. Fuhr, N. Kando, W. Kraaij, & A. de Vries (Eds.), *30th annual international ACM SIGIR conference* (pp. 527–534). ACM, July 2007. ISBN 987-1-59593-597-7.
- Stein, B., & Meyer zu Eissen, S. (2007). Intrinsic plagiarism analysis with meta learning. In B. Stein, M. Koppel, & E. Stamatatos (Eds.), *SIGIR workshop workshop on plagiarism analysis, authorship*

- identification, and near-duplicate detection (PAN 07)* (pp. 45–50). CEUR-WS.org, July 2007. URL <http://ceur-ws.org/Vol-276>.
- Stein, B., & Meyer zu Eissen, S. (2007). Topic-identifikation: Formalisierung, analyse und neue Verfahren. *KI—Künstliche Intelligenz*, 3, 16–22. ISSN 0933-1875. URL <http://www.kuenstliche-intelligenz.de/index.php?id=7758>.
- Stein, B., Lipka, N., & Meyer zu Eissen, S. (2008). Meta analysis within authorship verification. In A. M. Tjoa & R. R. Wagner (Eds.), *19th international conference on database and expert systems applications (DEXA 08)* (pp. 34–39). IEEE, September 2008. ISBN 978-0-7695-3299-8. doi: [10.1109/DEXA.2008.20](https://doi.org/10.1109/DEXA.2008.20).
- Surdulescu R. (2004). Verifying authorship. Final project report CS391L, University of Texas at Austin
- Tax, D. M. J. (2001). *One-class classification*. Ph.D. thesis, Technische Universiteit Delft.
- Tax D. M. J., & Duin, R. P. W. (2001). Combining one-class classifiers. In *Proceedings of the second international workshop on multiple classifier systems* (pp. 299–308). New York: Springer. ISBN 3-540-42284-6.
- Tweedie, F. J., & Baayen, H. R. (1998). How variable may a constant be? measures of lexical richness in perspective. *Computers and the Humanities* 32(5):323–352. doi:[10.1023/A:1001749303137](https://doi.org/10.1023/A:1001749303137).
- van Halteren, H. (2004). Linguistic profiling for author recognition and verification. In *ACL '04: Proceedings of the 42nd annual meeting on association for computational linguistics* (pp. 199). Morristown, NJ, USA: Association for Computational Linguistics. doi:[10.3115/1218955.1218981](https://doi.org/10.3115/1218955.1218981).
- van Halteren, H. (2007). Author verification by linguistic profiling: An exploration of the parameter space. *ACM Transactions on Speech and Language Processing*, 4(1), 1. ISSN 1550-4875. doi: [10.1145/1187415.1187416](https://doi.org/10.1145/1187415.1187416).
- Yang, H., & Callan, J. P. (2006). Near-duplicate detection by instance-level constrained clustering. In E. N. Efthimiadis, S. Dumais, D. Hawking, & K. Järvelin (Eds.), *SIGIR '06: Proceedings of the 29th annual international ACM SIGIR conference on research and development in information retrieval* (pp. 421–428). ISBN 1-59593-369-7.
- Yule, G. (1944). *The statistical study of literary vocabulary*. Cambridge: Cambridge University Press
- Zheng, R., Li, J., Chen, H., & Huang, Z. (2006). 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. doi:[10.1002/asi.20316](https://doi.org/10.1002/asi.20316).
- Zipf, G. K. (1932). *Selective studies and the principle of relative frequency in language*.