

13 description of the style of each author is in general subjective, and therefore
14 hard to derive in natural language; it is even harder to find a description
15 which enables a machine to automatically tell one author from the other. A
16 literature review on modern authorship attribution methods, usually coming
17 from the fields of machine learning and statistical analysis, is reported in
18 Stamatatos (2009); Jockers and Witten (2010); Koppel et al. (2009); Grieve
19 (2007); Juola (2006). Among these, algorithms based on similarity measures
20 such as Benedetto et al. (2002) and Koppel et al. (2011) are widely employed
21 and usually assign an anonymous text to the author of the most similar
22 document in the training data.

23 During the last decade, compression-based distance measures have been
24 effectively applied to cluster texts written by different authors (Cilibrasi and Vitányi,
25 2005) and to perform plagiarism detection (Chen et al., 2004). Such univer-
26 sal similarity measures, of which the most well-known is the Normalized
27 Compression Distance (NCD), employ general compressors to estimate the
28 amount of shared information between two objects. Similar concepts are
29 also used by methods using runlength histograms to retrieve and classify
30 documents (Gordo et al., 2013). Experiments carried out in Oliveira et al.
31 (2013) conclude that NCD-based methods for authorship analysis outper-
32 form state-of-the-art classification methodologies such as Support Vector
33 Machines. A study on larger and more statistically meaningful datasets
34 shows NCD-methods to be competitive with respect to the state of the art
35 (de Graaff, 2012), while Stamatatos (2009) reports that compression-based