

Quantitative patterns of stylistic influence in the evolution of literature

James M. Hughes^a, Nicholas J. Foti^a, David C. Krakauer^{b,c}, and Daniel N. Rockmore^{a,b,d,e,1}

^aDepartment of Computer Science, Dartmouth College, Hanover, NH 03755; ^bSanta Fe Institute, Santa Fe, NM 87501; ^cWisconsin Institute for Discovery, University of Wisconsin, Madison, WI 53715; ^dDepartment of Mathematics, Dartmouth College, Hanover, NH 03755; and ^eNeukom Institute for Computational Science, Dartmouth College, Hanover, NH 03755

Edited by* Michael S. Gazzaniga, University of California, Santa Barbara, Santa Barbara, CA, and approved March 13, 2012 (received for review September 21, 2011)

Literature is a form of expression whose temporal structure, both in content and style, provides a historical record of the evolution of culture. In this work we take on a quantitative analysis of literary style and conduct the first large-scale temporal stylometric study of literature by using the vast holdings in the Project Gutenberg Digital Library corpus. We find temporal stylistic localization among authors through the analysis of the similarity structure in feature vectors derived from content-free word usage, nonhomogeneous decay rates of stylistic influence, and an accelerating rate of decay of influence among modern authors. Within a given time period we also find evidence for stylistic coherence with a given literary topic, such that writers in different fields adopt different literary styles. This study gives quantitative support to the notion of a literary “style of a time” with a strong trend toward increasingly contemporaneous stylistic influence.

cultural evolution | stylometry | culture | complexity | big data

Written works, or literature, provide one of the great bodies of cultural artifacts. The analysis of literature typically involves the aggregation of information on several levels, ranging from words to sentences and even larger scale properties of temporal narratives such as structure, plot, and the use of irony and metaphor (1–3). Quantitative methods have long been applied to literature, most notably in the analysis of style, which can be traced back to a comment by the mathematician Augustus de Morgan regarding the attribution of the Pauline epistles (4) and the late nineteenth century work of the historian of philosophy Wincenty Lutasłowski, who brought basic statistical ideas of word usage to the problem of dating the dialogues of Plato (5). It was Lutasłowski who coined the word “stylometry” to describe such an approach to investigating questions of literary style. Since then, a wide range of statistical techniques have been developed toward this end (6), generally with the goal of settling questions of author attribution (see, e.g., refs. 6–11). Stylometric studies have also been pursued in the study of visual art (12, 13) and music [both in composition (14–16) and performance (17)], and are part of a growing body of work in the quantitative analysis of cultural artifacts (18).

In this paper we report our findings from the first large-scale stylometric analysis of literature. The goal of this work is not author attribution—for the authorship of all the works is well known—but is instead to articulate, in a quantitative fashion, large-scale temporal trends in literary (i.e., writing) style. This type of study has been, until now, impossible to undertake, but the advent of mass digitization has created dramatic new opportunities for scholarly studies in literature as well as in other disciplines (19). Our literature sample is obtained from the Project Gutenberg Digital Library (<http://www.gutenberg.org/wiki/Gutenberg>About>). Project Gutenberg consists of more than 30,000 public domain texts, music, audiobooks, etc., is freely available online, and is among the digital archives that have become, over the past 60 yr, crucial components of the preservation of cultural artifacts (18).

In scope, our work is related to, but quite different from, recent studies in the dating of literary works (20), the analysis of the coarse-grained structure of literary history (and the evolution of genre) (21), and most notably, a recent analysis of Google Books (22), wherein the temporal trends in content-word usage were articulated. Content words also form the basis of the various topic model analyses that have been applied to large corpora of science texts (see, e.g., ref. 23).

In contrast, the work presented here focuses on the usage of content-free words as the basis of the first large-scale study of the similarity structure of literary style. Content-free words are the “syntactic glue” of a language: They are words that carry little meaning on their own but form the bridge between words that convey meaning. Their joint frequency of usage is known to provide a useful stylistic fingerprint for authorship (8, 11), and thus suggests a method of comparing author styles. When we consider content-free word frequencies from a large number of authors and works over a long period of time, we can ask questions related to temporal trends in similarity. The primary results of our analysis are that time provides the most coherent means of clustering work and a trend of diminishing stylistic influence as we move forward in time. Such a finding is consistent with a simple evolutionary model for stylistic influence, which assumes that imitation attends preferentially to contemporary authors. In addition, we uncover quantitative support of the previously purely anecdotal notion of a literary “style of a time.” Taken together, these two findings suggest the utility and perhaps the creation of a new field of stylometric analysis in culturomics.

Materials and Methods

In our experiments, we studied a subset of the authors in the Project Gutenberg database composed of those who wrote after the year 1550, had at least five works in English in the Project Gutenberg collection, and for whom we had birth and death date information. This left us with 537 authors. For each author, we created a representative feature vector by aggregating the content-free word frequencies for each individual work by that author. In total, we analyzed 7,733 works.

In our experiments, we used a list of 307 content-free words that included prepositions, articles, conjunctions, “to be” verbs, and some common nouns and pronouns (see Table S1 for a complete list). We did not attempt to semantically disambiguate between occurrences of homographs in situ (e.g., when using “to” as a preposition or as an indicator of an infinitive verb). Doing so would require a sophisticated grammatical model, and it was not our aim to model this particular aspect of word usage. We believe that ignoring these distinctions is not likely to greatly affect our results, because words that account for the greatest differences in usage frequency among

Author contributions: J.M.H., N.J.F., D.C.K., and D.N.R. designed research; J.M.H., N.J.F., D.C.K., and D.N.R. performed research; J.M.H., N.J.F., D.C.K., and D.N.R. contributed new reagents/analytic tools; J.M.H., N.J.F., D.C.K., and D.N.R. analyzed data; and J.M.H., N.J.F., D.C.K., and D.N.R. wrote the paper.

The authors declare no conflict of interest.

*This Direct Submission article had a prearranged editor.

¹To whom correspondence should be addressed. E-mail: rockmore@cs.dartmouth.edu.

This article contains supporting information online at www.pnas.org/lookup/suppl/doi:10.1073/pnas.1115407109/-DCSupplemental.

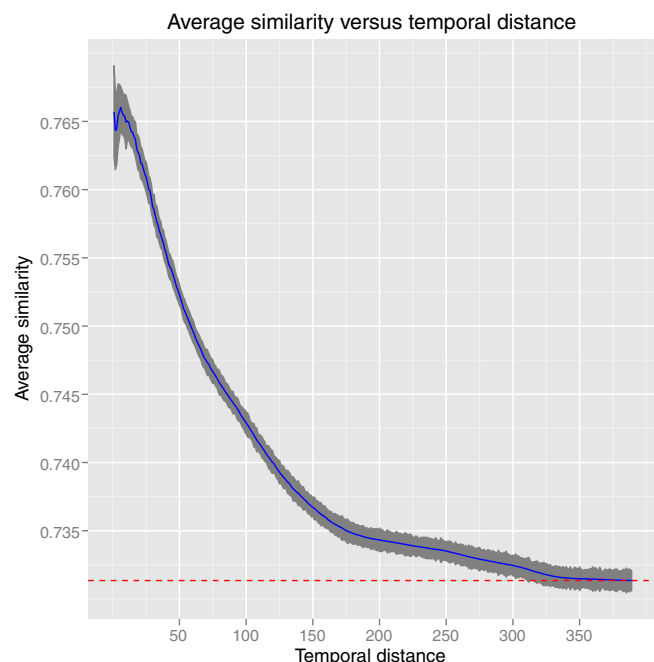


Fig. 2. Average similarity between authors as a function of temporal distance between them. Clearly, as the distance between authors increases, the similarity between them tends to decrease. The flat dashed red line marks the global average.

two time periods of equal length, “early” and “late.” For the early authors (those who wrote between 1550 and 1783), the average similarity as a function of temporal distance does not deviate significantly from the overall average, suggesting that authors during this time period influence each other roughly equally, regardless of how far apart in time they are.

However, for the “modern” authors (those who wrote between 1784 and 1952), the average similarity curve was high for shorter temporal distances and decreased rapidly toward the mean, much like the overall trend shown in Fig. 2. In order to examine the extent to which there is a shift in the way authors were influenced based on when they wrote, we split the modern authors into quartiles defined over the range of most densely populated years (see Fig. 3) by partitioning them according to their representative author years. The four partitions consist of the years 1784–1829, 1825–1870, 1866–1911, and 1907–1952. There is a small amount of overlap (5 yr) between these groups in order to mitigate edge effects.

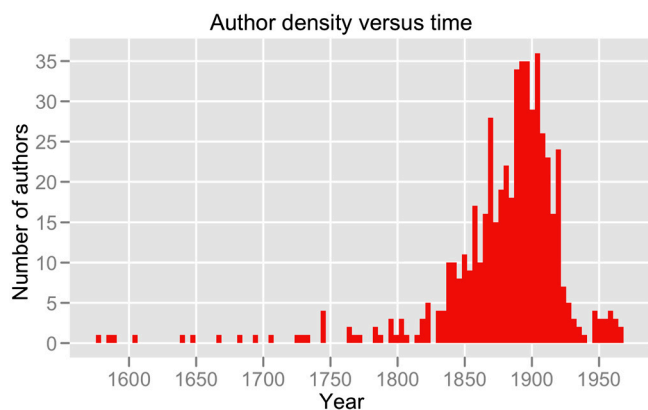


Fig. 3. Density of authors in our dataset as a function of time. The vertical axis indicates how many authors fell into the corresponding time window.

Fig. 4 *A* and *B* display $S_{\text{avg}}(t)$ and $S_W(t)$ for each of these groups separately, indicating that they possess remarkably different patterns of similarity as a function of temporal distance.

In Fig. 4*A* we plot $S_{\text{avg}}(t)$ for all authors in the corresponding time window. For the early modern period (1784–1870), the similarity functions do not differ significantly from the average (indicated by the dashed red line), which suggests that authors during that period tended to draw influence from other authors uniformly as a function of temporal distance. The same pattern is observed for the windowed analysis in the same time period (see Fig. 4*B*). Thus over this period there is no significant evidence for stylistic localization in time.

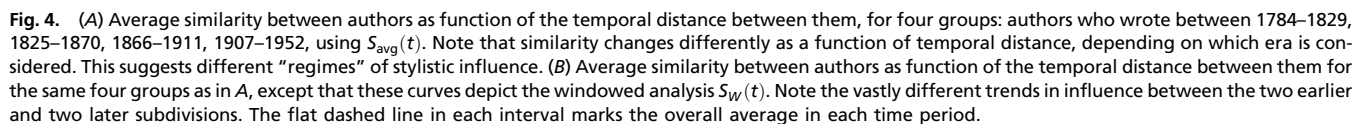
For the late modern quartiles (1866–1911 and 1907–1952) the pattern is very different. In the period 1866–1911, authors are significantly more similar to members of their own age cohort, and interestingly this similarity decreases toward the average, for the cumulative analysis. In other words, above 30 yr apart, authors are not significantly more like any set of authors chosen at random, regardless of how far apart they are in time. When we consider the windowed analysis, we see more structure. Above 20 yr apart, authors actually tend to be less like each other than the average. This suggests that there is a significant decline in the similarity between authors who are widely separated in time. The “repulsive” effect of temporal distance is consistent within this period, and similarity between authors decreases throughout the years in question.

In the later period, 1907–1952, this pattern repeats but with stronger effects. Contemporaneous authors are most similar and average similarity decays to the within-group average with increasing temporal separation. The rate of decay is now nonlinear and scales quadratically in time ($s \approx t^{-2}$), which suggests that authors tend to be influenced by their contemporaries more strongly than during 1866–1911. In other words, the amount of temporal distance until contributions of authors becomes indistinguishable from the average is smaller than in the earlier period spanning 1866–1911. In the later period, average similarity was indistinguishable from the mean after approximately 23 yr, whereas in the earlier period, it was not until almost 30 yr that influence became random. The windowed analysis, as before, exhibits greater structure. Average similarity $S_W(t)$ is no longer monotonic in time but has a minimum at 25 yr, returning to the average at the maximum calculated separation of 40 yr. This suggests that for the modern period the pattern observed over the complete dataset is reversed. Whereas when we consider the average similarity over the complete dataset, the most similar authors are around 24 yr apart, in the late modern quartile, authors separated by 25 yr are maximally different. These findings are robust under changes in sampling (see Fig. S1).

Discussion

It is a remarkable fact that vectors of content-free words—subject-independent textual features of a book—allow us to cluster authors in time and by narrative theme, and that content-free word frequencies are fairly faithfully transmitted among authors of a similar period, even when imitation at this level of textual resolution seems to be out of the question. As we move into the present, this imitation becomes increasingly localized to our contemporaries.

We propose that for the earliest periods in our dataset, and the early modern period, the number of published works remained relatively low. This allowed authors to have sufficient time to sample (read) very broadly from the full range of historically published works. Common phrasing, and norms of syntax and grammar, remain relatively unchanged for long periods of time. This generates decay rates in similarity as a function of temporal distance that are not significantly different from the average, because authors are influenced by models distributed uniformly in time. However, for more recent authors, the number of possible



The patterns of influence are a first discovery from the corpus. Implicit in this is a temporal clustering of similarity and quantitative support for the qualitative suggestions of a notion of a “style of a time.” It is also worth noting that the implicit temporal clustering of similarity is not an exclusively temporal phenomenon. Fig. S2 shows a network representation of the authors in which a preliminary investigation reveals evidence of thematic clustering as well. Examples include interesting groupings of

ACKNOWLEDGMENTS. D.C.K.'s contribution to this project/publication was made possible through the support of a grant from the John Templeton Foundation.

- PNAS | May 15, 2012 | vol. 109 | no. 20 | 7685

