

Text Analysis: A Survey of Principles, Tools, and New Ways of Reading

Christopher Ohge
School of Advanced Study, University of London

Email: christopher.ohge@sas.ac.uk
Twitter: @cmohge

What is text analysis?

- 1. Using computational tools and/or programs to analyse text data of varying sizes by yielding quantitative results, such as simple character and word frequencies, mean word usage, and type-token and hapax percentages (to identify repetition, brevity, and unique word clusters.)*
- 2. Making arguments about historical trends or form, style, sentiment, and/or context.*

This can be broken down further:

Statistical analysis: counting relative word frequencies, calculating word pairings (collocates) and ngrams (groups of two or three words or more co-occurring words), average word and sentence length, mean word usage; organising words that occur too frequently to be studied one by one.

Corpus analysis: searching and pattern recognition across multiple texts.

Linguistic analysis: parts-of-speech tagging, lexical variety and uniqueness, topic modeling, sentiment analysis, stylometry.

Network analysis: mapping connections between metadata and textual data.

Various visualisations (graphs, maps, and trees) for explanatory force.

Some principles

- 1. Relative frequency:** cf. John Stuart Mill's **principle of concomitant variation**, which states that the magnitude of an antecedent variable changes in proportion to the magnitude of a second variable. Effects are typically proportional to their causes. This means that texts with similar relative frequency ratios tend to be correlated and may be causally connected (e.g., authored by the same person).
- 2. Exploratory data analysis (cf. John Tukey):** clustering similar groups of results based on simple statistics and geographical representations (using, e.g., Burrows Delta, Cluster Analysis, Principal Component Analysis).
- 3. Inductive inference and probability:** creating general hypotheses from specific instances. Yet, understanding that the results are only probabilistic (i.e., only as good as the quality and comprehensiveness of the data). Justified true belief rather than mere fact.
- 4. Computational close reading.** Looking closer at peripheral results, disaffinity and exceptions; letting the research questions guide the exploration; and exercising skepticism.
- 5. Distant reading:** Making broad claims about historical text data by analysing samples of huge data sets.

What is the point of studying the most frequent common words?

The most frequent words in Herman Melville's *Moby-Dick*, visualised in Voyant Tools (voyant-tools.org)



		Term	Count
⊕	□	1 the	144...
⊕	□	2 of	6609
⊕	□	3 and	6430
⊕	□	4 a	4715
⊕	□	5 to	4625
⊕	□	6 in	4172
⊕	□	7 that	2990
⊕	□	8 his	2530
⊕	□	9 it	2420
⊕	□	10 i	1989

The value of studying common words

Anthony Kenny discusses the “mysterious veneration” that some literary scholars have for single and rare word occurrences, when “the rate of occurrence of a dull common word in a text may be a much more significant feature” (*Computation of Style* [1982], 67–68).

Similarly, **John F. Burrows**, in *Computation into Criticism* (1987), bases his analysis on the 30 most common words in Jane Austen’s novels, with less attention to unique words.

Burrows (2004) later argues that the styles of authors come from common words (i. e., articles and prepositions).

“[T]he real value of studying the common words rests on the fact that they constitute the underlying fabric of a text, a barely visible web that gives shape to whatever is being said ... The principal point of interest is neither a single stitch, a single thread, nor even a single color but the overall effect. Such effects are best seen, moreover, when different pieces are put side by side.”

Burrows: “computer-assisted textual analysis can be of value in many different sorts of literary inquiry, helping to resolve some questions, to carry others forward, and to open entirely new ones.”

Distant reading and interpretation

Franco Moretti: “Quantitative research provides a type which is ideally independent of interpretations … it provides *data*, not interpretation” (*Graphs, Maps, Tress* [Verso, 2007]).

History: shifting the gaze from extraordinary people and events to everyday facts. What literature can be found in large mass of facts?

“Abstraction is not an end in itself, but a way to widen the domain of the literary historian, and enrich its internal problematic” (Moretti, p. 2).

Distance is a new kind of knowledge—a model that raises the level of abstraction to increase cognitive capacity.

Possible directions of travel

- Testing a hypothesis or thesis about an author, text, passage, genre, or period
- Testing the claims of a critical work
- Investigating how and the extent to which authors differentiate the voices of characters or narrators in a work
- Investigating shifts in style and how they change over time
- Investigating the history of an important word, concept, or group of words or concepts over a long time span
- Studying the effects of genre conventions
- Investigating claims of authorship

David Hoover (2013): “the computer’s greatest strengths are in storing, counting, comparing, sorting, and performing statistical analysis. This makes computer-assisted textual analysis especially appropriate and effective for investigating textual differences and similarities” (“Textual Analysis”).

Computer-assisted versus computational

Computer-assisted: using out-of-the-box tools to generate results.

Pros: fast results, good interfaces and visualisations.

Cons: lack of control over features and inability to manipulate data; hard to validate and replicate results.

Computational: using programming languages to generate results.

Pros: you have as much control as you can muster; complete customisation.

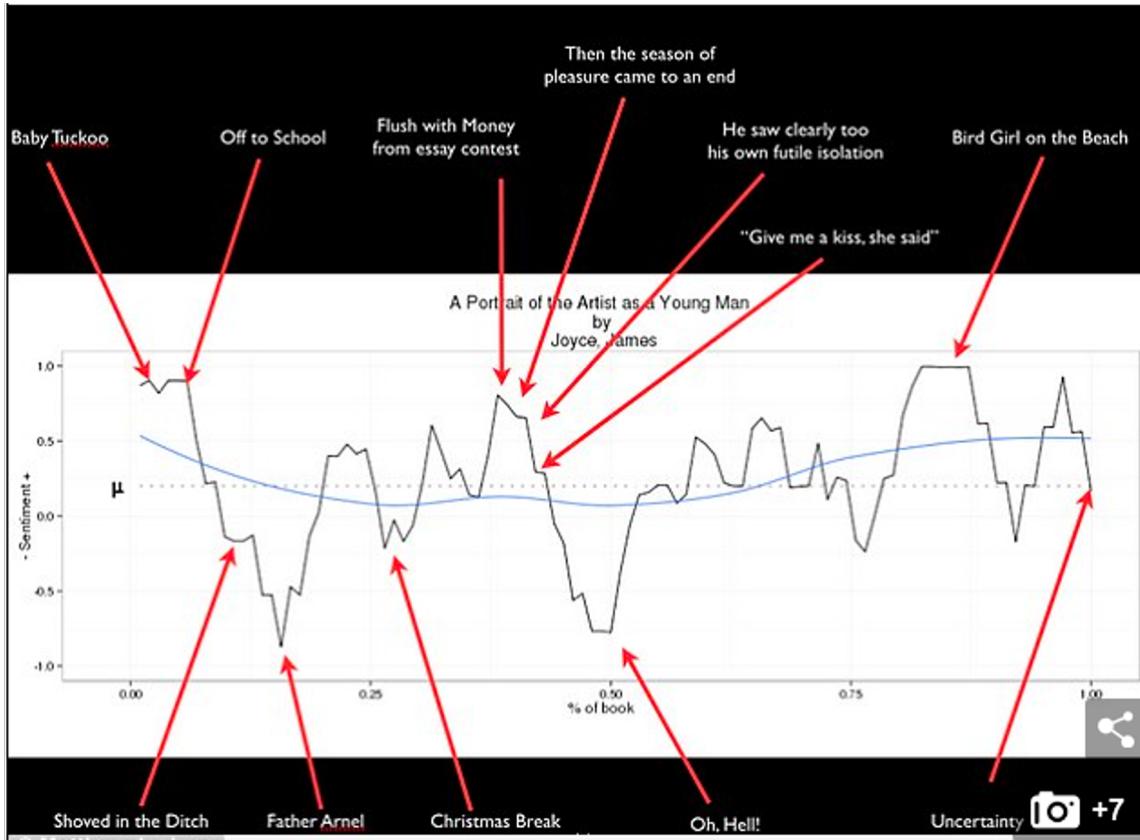
Cons: requires some knowledge of programming.

Options for performing text analyses

- Voyant Tools <<https://voyant-tools.org/>>
- AntConc <<http://www.laurenceanthony.net/software/antconc/>>
- Language databases, such as the Historical Thesaurus of English
[<https://ht.ac.uk/>.](https://ht.ac.uk/)
- Text database tools, such as Hathi Trust Bookworm and Google nGram searches.
- Programming Languages: R <<https://www.r-project.org/>> or Python
 [<https://www.python.org/](https://www.python.org/)

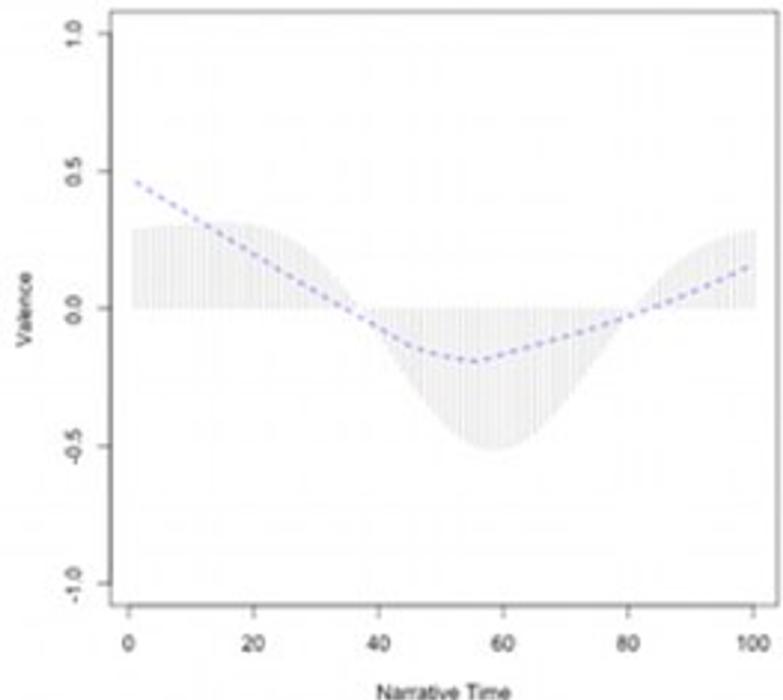
Distant Reading:

'Professor who analysed 40,000 novels claims there are just SIX possible storylines' (*Daily Mail*, 26 February 2015).



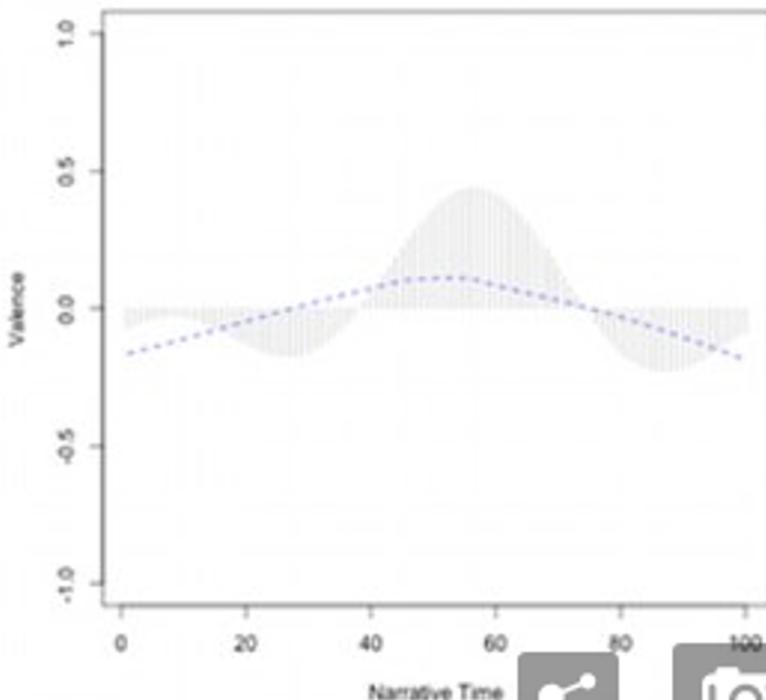
“Man in Hole”

19031 Books: (Mean Shape for 45.99% of Corpus)

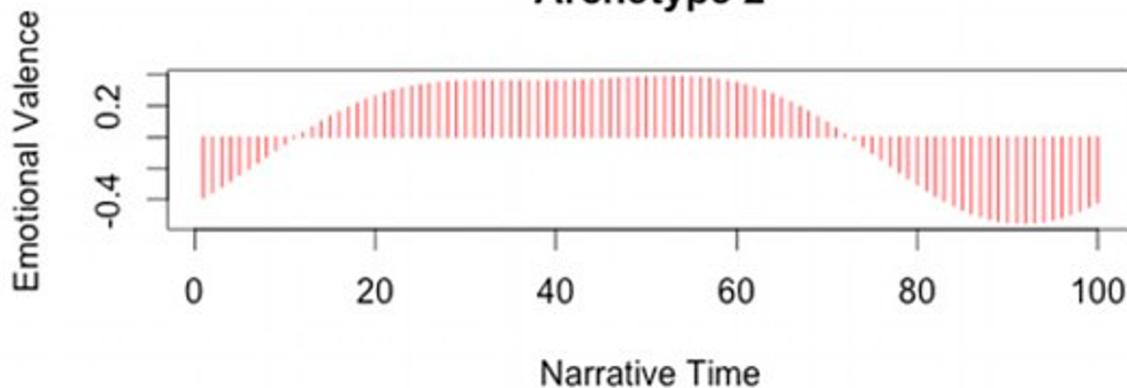


“Man on Hill”

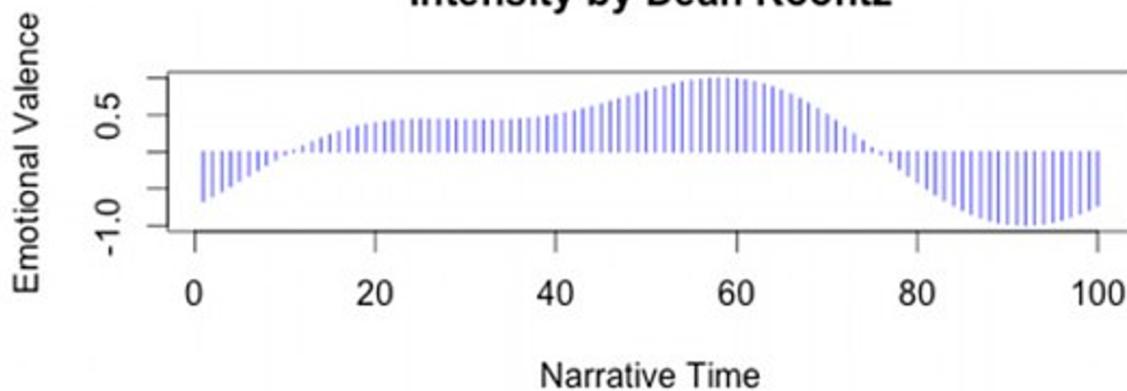
22352 Books: (Mean Shape for 54.01% of Corpus)

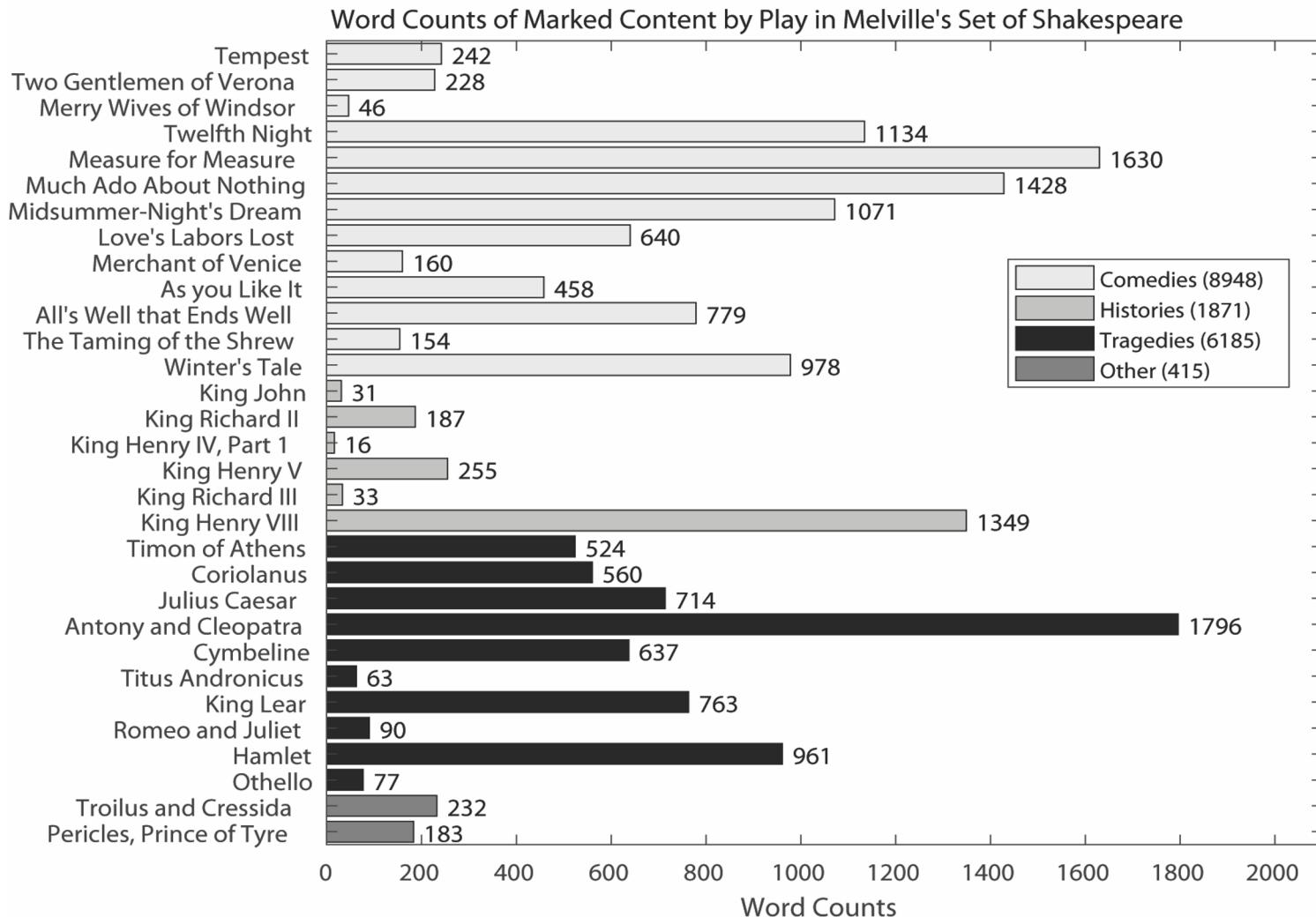


Archetype 2



Intensity by Dean Koontz





Lexical Uniqueness Values for each Marked Passage in Melville's Marginalia to Timon of Athens, King Lear, Hamlet, and Othello

