



HAL
open science

An Experiment on Plato's Gorgias as an Introduction to Textometry

Bénédicte Pincemin, Stéphane Marchand

► **To cite this version:**

Bénédicte Pincemin, Stéphane Marchand. An Experiment on Plato's Gorgias as an Introduction to Textometry. *Classics@*, 2022, Digital Text Analysis of Greek and Latin sources, 20. halshs-01730373

HAL Id: halshs-01730373

<https://shs.hal.science/halshs-01730373>

Submitted on 26 Aug 2022

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

An Experiment on Plato's *Gorgias* as an Introduction to Textometry

Bénédicte Pincemin, Stéphane Marchand

1. Background: What is Textometry? How could it concern Classical Studies?

Current implementations of Digital Humanities for classical studies are often either rich digital editions, with much effort devoted to display a pleasant layout and to offer efficient navigation paths; or indexed corpora or databases with querying facilities, for which the focus is on searching functionalities and not on primary text visualization and reading. However, through the example of the textometric approach (presented in section 1.1) that is implemented in TXM open-source software (1.2), we would like to show that complex textual representation and computational analysis can be both taken into consideration into a unique digital framework (1.3). In order to show that such a method can offer great possibilities for researchers in the Classics field, we will give a more complete presentation of the methodology applied to Plato's *Gorgias* (section 2 and 3).

1.1 What is textometry?

Textometry is a computer-assisted textual analysis method that is based on word counts (or on any linguistic feature counts) and that links statistical processing to data referencing and contextual comparison. Main textometric functions are specificities (a keyword computation), correspondence analysis (a geometrical view of the corpus content), cooccurrences (lexical collocations), and KWIC (Keyword in context) concordances (for a technical presentation of these functions, see appendix, section 5).

The textometric methodology was founded in the 70s in France (Léon and Loiseau 2016). Benzécri and his students launched the *analyse des données* field and created new exploratory

multivariate statistical methods, mainly correspondence analysis and its combination with clustering (Benzécri 1973, 1981). Then, the Saint-Cloud Laboratory coined the word “lexicometry” (which would evolve later to “textometry” and “logometry”, but the core procedures remain the same) and added new methods especially dedicated to quantitative lexical analysis: specificities (*spécificités*), collocations (*cooccurrences*) (Lafon 1984), phrase detection with repeated segments (*segments répétés*) (Salem 1987). At Nice University, Brunet was also a pioneer in developing and applying these new methods to French literature (Brunet 2011). An overview of the methods is gathered in Lebart, Salem, and Berry 1998. Moreover, since 1992, the JADT Conference (International conference on statistical analysis of textual data), which is led by the textometric scientific community, has been held every two years; its proceedings have been published online since 1998 on the Lexicométrica website.¹

Compared to other text analysis methods, textometry is a balanced combination of quantitative and statistic methods on the one hand, and qualitative methods on the other hand. It occupies an intermediate position between text mining (which replaces the textual data by quantitative summaries, extractions, and visualizations) and annotation software (which enrich and investigate a detailed view of the text); in other words it achieves both a close reading and a distant reading.

Several textometric applications are currently available for research, like DtmVic, Hyperbase, Lexico, IRaMuTeQ, Trameur, TXM.² In this paper our example will be run with TXM.

1.2 What is TXM?

TXM is an open-source textometric software that was initiated in 2009 within the frame of the Textométrie project (2007–2010, founded by the ANR French agency), gathering four research laboratories, from Lyon, Paris, Nice, and Besançon. Its development is coordinated by Serge Heiden (Lyon, IHRIM laboratory), and the two main developers are Matthieu Decorde (Lyon, IHRIM Laboratory) and Sébastien Jacquot (Besançon, ELLIAD Laboratory) (Heiden 2010).³

¹ JADT conference archives: <http://lexicométrica.univ-paris3.fr/jadt/index.htm>

² Here are the websites for these applications: DtmVic: <http://www.dtmvic.com> ; Hyperbase: <http://ancilla.unice.fr/hyperbase> ; Hyperbase Web Edition: <http://hyperbase.unice.fr> ; IRaMuTeQ: <http://www.iramuteq.org> ; Lexico 5: <http://www.lexi-co.com> ; Trameur: <http://www.tal.univ-paris3.fr/trameur> ; TXM: see next footnote.

³ Textometry project and TXM software website: <http://textometrie.org>

The aim of the Textometry project was for the new software to be able to fully manage and analyze state-of-art corpora, that is, structured and annotated corpora, like TEI corpora. Another aim was to launch an open-source development the partners could share.

TXM implements the textometric approach, so that it includes specificities, correspondence analysis, cooccurrences, and back-to-text facilities implemented as an advanced concordance view linked to a text view. This paper will show some other available functions too: progression in section 3.2, and index (word frequency lists) mainly in 3.4. Indexing and searching are operated by the CQP⁴ component (Christ 1994; Evert and Hardie 2011), and statistical computing relies on R⁵ components; processing can be extended or personalized with scripts. One of TXM main features is its ability to input many kinds of digital texts, from plain text to generic XML and XML-TEI. On the basis of these features, TXM is a fifth generation concordancer according to McEnery and Hardie typology (McEnery and Hardie 2012).

TXM runs on Windows, Linux, MacOS, and can also be set as a web portal version, for collaborative work or on-line corpus publication.⁶

1.3 Three (Functions) in One (Digital Edition): Reading, Searching, Analyzing

As textometry pays attention both to complex queries and statistics, and to fine text reading, it can be a new unified framework for digital editions in a classical studies context. With a TXM solution, for example, humanists are not compelled to choose between either a digital edition, or a database, or a search engine, to study and publish their text corpora.

Actually, TXM is a *search engine*. Complex queries can be processed about words, phrases, morphosyntactic patterns (or patterns based on any other information that has been added to the corpus through tagging), also allowing for search patterns that include gaps (like, for instance, in ancient Greek μὲν ... δὲ). Queries can cross information from every level, from lexical tags to text structures, which can be internal structures (inside texts, like chapters) or overall structures about groups of texts (like author or text genre).

⁴ CQP stands for Corpus Query Processor, which is the name of the search engine component included in the IMS Open Corpus Workbench (CWB) (see <http://cwb.sourceforge.net>) initially developed at the University of Stuttgart (<http://www.ims.uni-stuttgart.de/forschung/projekte/CorpusWorkbench>).

⁵ R is a freely available language and environment for statistical computing and graphics. CRAN (The Comprehensive R Archive Network, <http://cran.r-project.org>) is currently one of the main providers for open-source and up-to-date statistical packages.

⁶ The developer environment and some of user documentation (user manual [Heiden, Decorde, and Jacquot 2018], user interface, mailing list) are available in English.

At the same time TXM is a *database* that manages text documents. Numerous metadata can be associated to each document (such as author, date of composition, text type, etc.), and selection of documents can be built based on the metadata's values. This can be illustrated by the text selection function in the Base de français médiéval instance of TXM web portal version (Figure 1),⁷ which shows very detailed selection panels that are used like database query forms. Furthermore, any textual unit, from words or text chunks to text groups, can be described with metadata information, and then be used to target a selection in combination with any information on any text level. For instance, one may focus on sentences which include some linguistic pattern and which were written within a given time span.

The screenshot displays the TXM web portal interface for the Base de français médiéval. The main window shows a search results table with columns for 'notice', 'T', 'auteur', 'titre', 'date compo libre', 'ms date li...', 'forme', and 'domaine'. The table lists various medieval texts, such as 'Serments de Strasbourg', 'Séquence de sainte Eulalie', and 'Chanson de Roland'. On the left side, there are several selection panels for metadata fields like 'auteur', 'titre', 'siècle', 'date oeuvre', 'forme', 'domaine', 'genre', 'dialecte', 'relation', 'morphosynt', and 'restrictions'. Each panel has a 'valeur' field and a 'n' field for the number of results. The 'forme' panel shows a table with columns 'valeur', 'n', 'N', 't', and 'T', with rows for 'mixte', 'prose', and 'vers'. The 'domaine' panel shows a table with columns 'valeur', 'n', 'N', 't', and 'T', with rows for 'littéraire', 'politique', and 'religieux'. The 'genre' panel shows a table with columns 'valeur', 'n', 'N', 't', and 'T', with rows for 'romain', 'hagiographie', 'miracle', 'épopée', 'chronique', 'récit bref', 'bestiaire', 'lyrique', and 'sermon'. The 'dialecte' panel shows a table with columns 'valeur', 'n', 'N', 't', and 'T', with rows for 'champenois', 'nglo-normand', 'champenois méridional', 'normand', 'poitevin', 'picard', 'Ouest', 'Nord-Ouest', 'franco-picard', and 'lillois'. The 'relation' panel shows a table with columns 'valeur', 'n', 'N', 't', and 'T', with rows for 'relation' and 'morphosynt'. The 'restrictions' panel shows a table with columns 'valeur', 'n', 'N', 't', and 'T', with rows for 'restrictions' and 'morphosynt'. The table on the right shows the results of the search, with columns for 'notice', 'T', 'auteur', 'titre', 'date compo libre', 'ms date li...', 'forme', and 'domaine'. The table lists various medieval texts, such as 'Serments de Strasbourg', 'Séquence de sainte Eulalie', and 'Chanson de Roland'. The table is sorted by 'forme' and 'domaine'. The 'forme' column shows values like 'prose', 'vers', 'littéraire', and 'didactique'. The 'domaine' column shows values like 'juridique', 'religieux', 'littéraire', and 'didactique'. The table also shows the number of results for each row, with a 'n' column. The table is displayed in a grid format with checkboxes for each row. The interface also includes a search bar at the top, a navigation menu on the left, and a footer with the URL 'http://bfm.ens-lyon.fr'.

Figure 1. Text selection in the Base de Français Médiéval TXM web portal.

Then, TXM provides a wide range of corpus linguistic functions to process these lexical and textual selections (word frequency lists, KWIC concordance, collocations, keyword analysis, distributions, multidimensional analysis, etc.; see section 1.2). These analytic features can be coupled with a fine text encoding (like XML-TEI), making it possible to record precise philologic

⁷ Base de français médiéval: <http://bfm.ens-lyon.fr>

information⁸ and to build fine-tuned HTML *editions* to visualize the text (Figure 2). One important feature of the process is the possibility to manage and distinguish information used as text content (words to be searched and counted), and information useful to follow the history of the text, assess its structure, and assist text reading and interpretation (critical apparatus, editorial text divisions, speech turn labels, etc.). If available, facsimiles of source documents can be displayed so that the researcher still has a view on the primary document. Thus one keeps access to information that might have been lost because of necessary digital encoding choices.

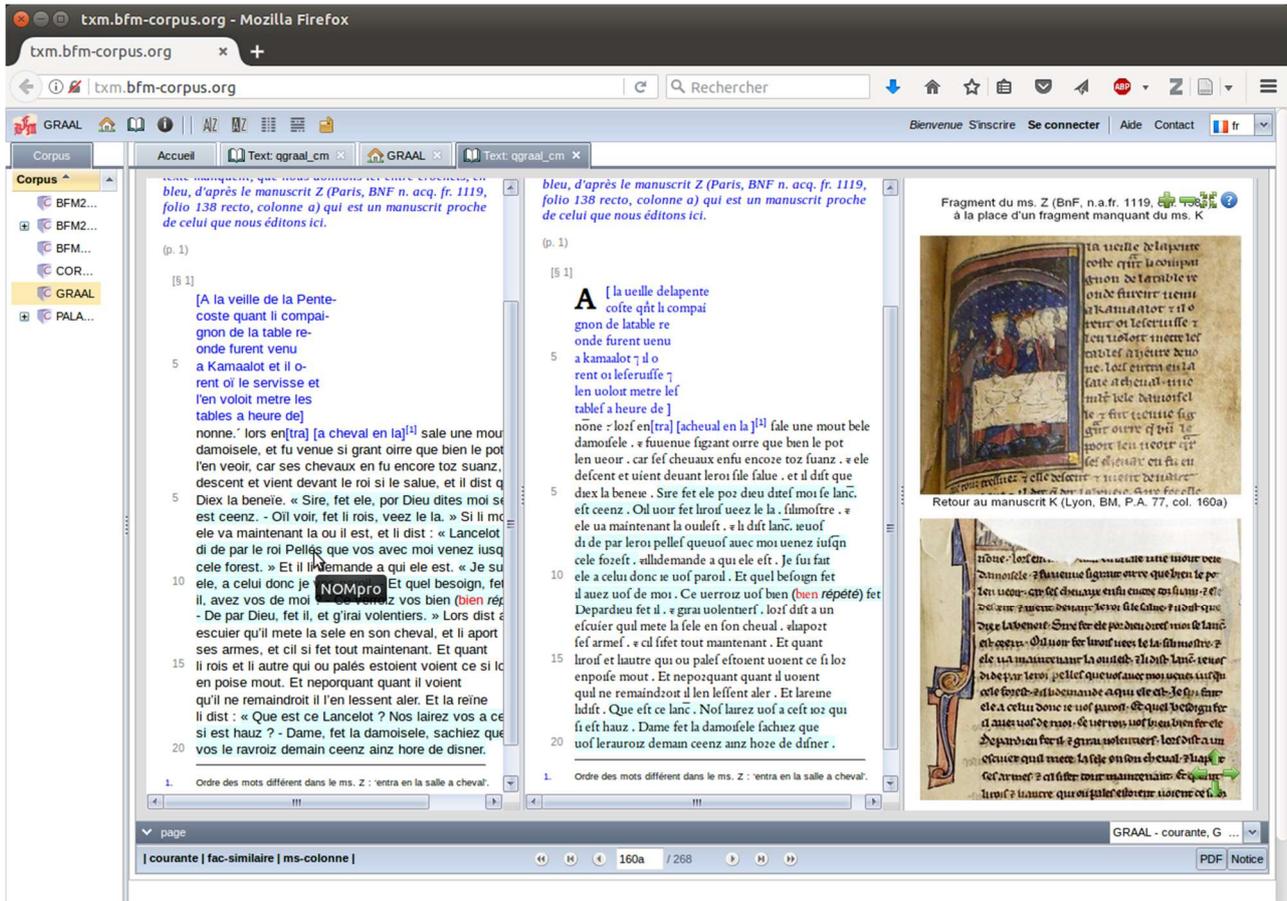


Figure 2. Synoptic views of the *Queste del Saint Graal* in a TXM web portal edition (Marchello-Nizia and Lavrentiev 2013): here two kinds of transcriptions and the manuscript image are aligned.

Thus with this example we see that it is possible to give the same text a unified digital edition that offers all together searching facilities, database querying, and an elaborated and precise layout. The aim of the two following sections is to present in concrete terms what kind of investigations may be done within such a framework, in the context of ancient Greek studies and current digital resources.

⁸ For instance, the TEI provides means to encode gaps, additions, deletions, notes, emphasis, etc.

2. Corpus Example: Plato's Work, from Perseus to TXM⁹

2.1 Choice of the texts

As a first experiment of TXM on ancient Greek literature, we have chosen to study a corpus of Plato's texts, because one of us has a good knowledge of these texts, especially of *Gorgias* (Marchand and Ponchon 2016). Actually, the textometric approach presupposes some acquaintance with the corpus to be investigated, because this approach does not compute an analysis on behalf of the researcher: one cannot give the corpus as input, run a (black box) program, and get what would be the complete textometric analysis for the corpus. On the contrary, it is the responsibility of the researcher to lead the investigation and to raise relevant questions so that the produced results can really be meaningful and fruitful. Since *Gorgias* is a rather small, well studied text, we already have some concrete ideas about its features before starting the textometric analysis, so we can check the results of our textometric analysis against our prior knowledge. Furthermore, the software allows one to perform queries that are impossible for a human reader, or at least that would consume a large amount of time. At any rate, we could take this software as an empirical tool that allows us to test some hypotheses, or to give us some clues to continue investigating such an intriguing dialogue as *Gorgias*, and Plato's work more widely.

We downloaded the XML-TEI edition of the Burnet edition of Plato's texts, published by the Perseus Digital Library¹⁰ and provided under a Creative Commons Share Alike 3.0 license. For this experiment, among the 36 available texts, we excluded the works that are generally considered as *spuria* or *dubia* and, for stylistic reason, the letters.¹¹ Hence the corpus is composed of 29 texts: *Euthyphro*, *Apology*, *Crito*, *Phaedo*, *Cratylus*, *Theaetetus*, *Sophist*, *Statesman*, *Parmenides*, *Philebus*, *Symposium*, *Phaedrus*, *Alcibiades 1*, *Alcibiades 2*, *Charmides*, *Laches*, *Lysis*, *Euthydemus*, *Protagoras*, *Gorgias*, *Meno*, *Hippias Major*, *Hippias Minor*, *Ion*, *Menexenus*, *Republic*, *Timaeus*, *Critias*, and *Laws*.

⁹ The import process followed here and the resources associated with it are available on txm-users wiki: <https://groupes.renater.fr/wiki/txm-users/public/perseus>. The page which is dedicated to this article (https://groupes.renater.fr/wiki/txm-users/public/perseus_201707_plato) provides also the binary (.txm) version of Plato's corpus, which can be directly loaded into TXM: this binary version allows one to skip the technical stages of the import process.

¹⁰ Perseus Home: <http://www.perseus.tufts.edu/hopper>; Text repository: <https://github.com/PerseusDL/canonical-greekLit>.

¹¹ Most of the *Letters* are considered as spurious. Moreover the *Letters* belong to a different text type and would disturb contrastive analysis within the corpus (we could study them separately, in another corpus).

2.2 From Perseus TEI encoding to TXM

When we prepared this corpus in June and July 2017, TEI encoding of Plato's texts in Perseus was heterogeneous. We had to deal with texts in several states: last updates made in 2017, 2015, 2014, and 1992. 2017 texts were clearly a new encoding generation. Two texts (*Ion* and *Republic*) used a particular encoding scheme, for instance regarding the use of the <div> element and the marking of the sections' beginning. We decided not to modify sources (which are continuously evolving and improving thanks to the Perseus community), but to make automatized and limited changes during the import process so as to get a usable corpus, even if the TXM user had to contend with some inherited heterogeneity (this is especially the case for speaker information). These automatic changes were processed with XSL and CSS stylesheets,¹² used as parameters in the XML-XTZ¹³ import.

The main stylesheet, applied in the second stage of the import process, manages some XML-TEI features of Perseus texts (about nested <div> or <text>) in order to make them compliant with TXM processing (especially for the CQP search engine component embedded in TXM). It automatically gets information from <teiHeader> and associates it with each text instance in TXM: <title>, <author>, and <editor> from <titleStmt> (first mention for each element), and also the content of @when attribute for first (or most recent) <change> element in <revisionDesc>.¹⁴

Another XSL stylesheet builds the default references provided for word matches in KWIC concordance view: we chose to show the usual name for the text followed by the Estienne page number (Figure 3). Additionally, we ensured that CTS-URN information, providing unique and standard identifiers to cite digital textual data, is also available at the word level, and can alternatively be chosen for references if needed (Figure 4). The main stylesheet was also used to add page break <pb> elements (with the page number) so as to have the same pages in TXM as in the Perseus edition.

¹² We are very grateful to Alexei Lavrentiev (CNRS, IHRIM, Lyon), who helped us to adapt the stylesheets to our corpus and needs.

¹³ XTZ stands for XML-TEI Zero, that is, this import is intended to process XML-TEI files, taking into account the semantics of some common elements (<text>, <p>, <ab>, <lg>, <head>, <lb>, <pb>, <w>, <graphic>, <ref>, <note>, <hi>, <emph>, <list>, <item>, <table>, <row>, <cell>) if they occur (no mandatory element except <text>). By default, other elements are kept for the analysis and have no rendering in corpus edition view, but this can be fine-tuned by means of XSL and CSS parameter stylesheets. More information is available online in TXM documentation: developer specifications (https://groupes.renater.fr/wiki/txm-info/public/import_xtz), user manual (<https://txm.gitpages.huma-num.fr/txm-manual/importer-un-corpus-dans-txm.html#module-xml-tei-zero-csv-dit-aussi-xtzcsv-ou-xtz-import-de-xml-tei-generique-.xml>).

¹⁴ A normalized form for short title (title1) and date (update10) is also provided by a metadata.csv parameter file.

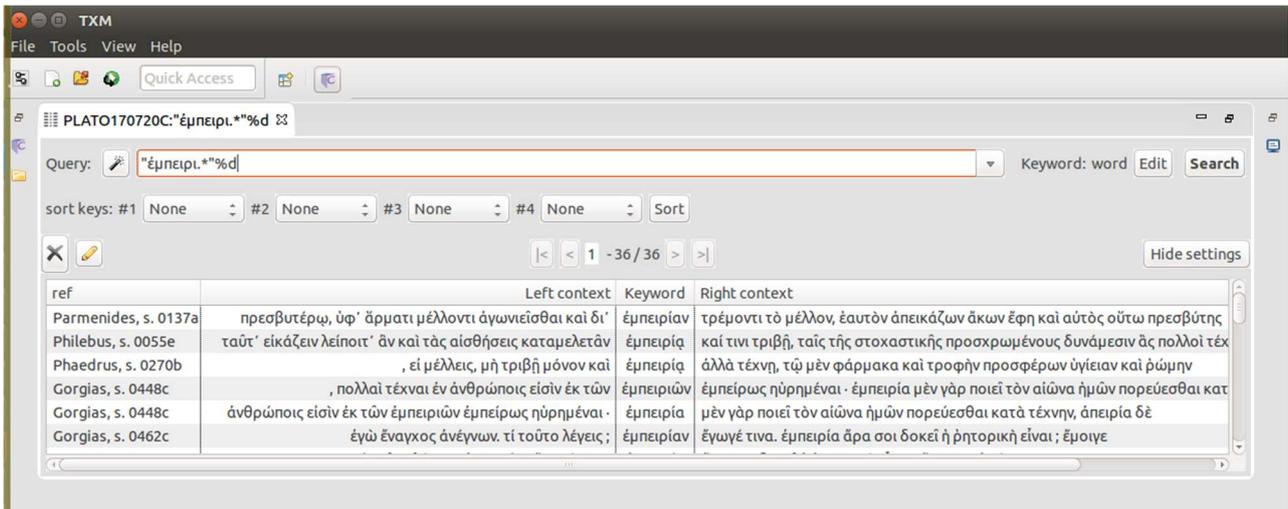


Figure 3. Concordance view with default references in the first column (short title followed by Estienne page number).¹⁵

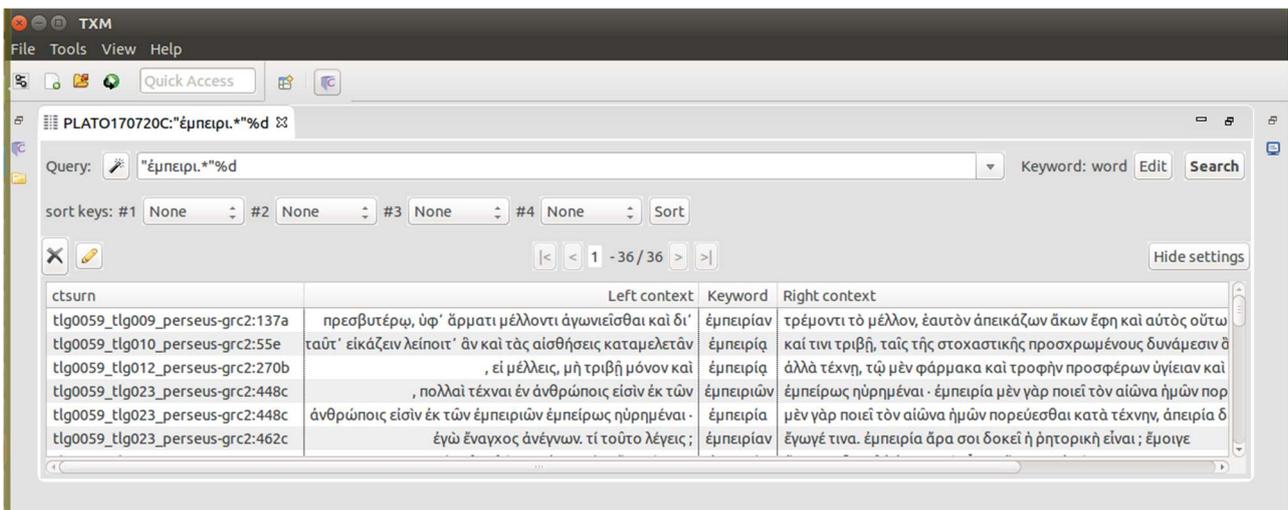


Figure 4. Concordance view with references set to CTS-URN (first column).

We wanted to keep and show clearly speech turn information and speaker information, without indexing the protagonist's name as a word to be counted and searched as such. To this aim we declared <speaker> and <label> elements as “Out-of-text-to-edit” element in import parameters. As text encoding in Perseus was heterogeneous, speaker information is also heterogeneous in TXM. Speech turns are defined by either <sp> or <said> element (depending on the text), and speaker is indicated with @who attribute only in <said> elements (but not all of them) for analysis purposes, and is displayed using <label> or <speaker> element content when available (Figure 5).

¹⁵ The query is: ""ἐμπειρι.*%d". It matches any word beginning with ἐμπειρι, and the %d operator allows diacritic variations, so that this query matches ἐμπειριῶν but also ἐμπειρίᾳ with an accent on the last ι.

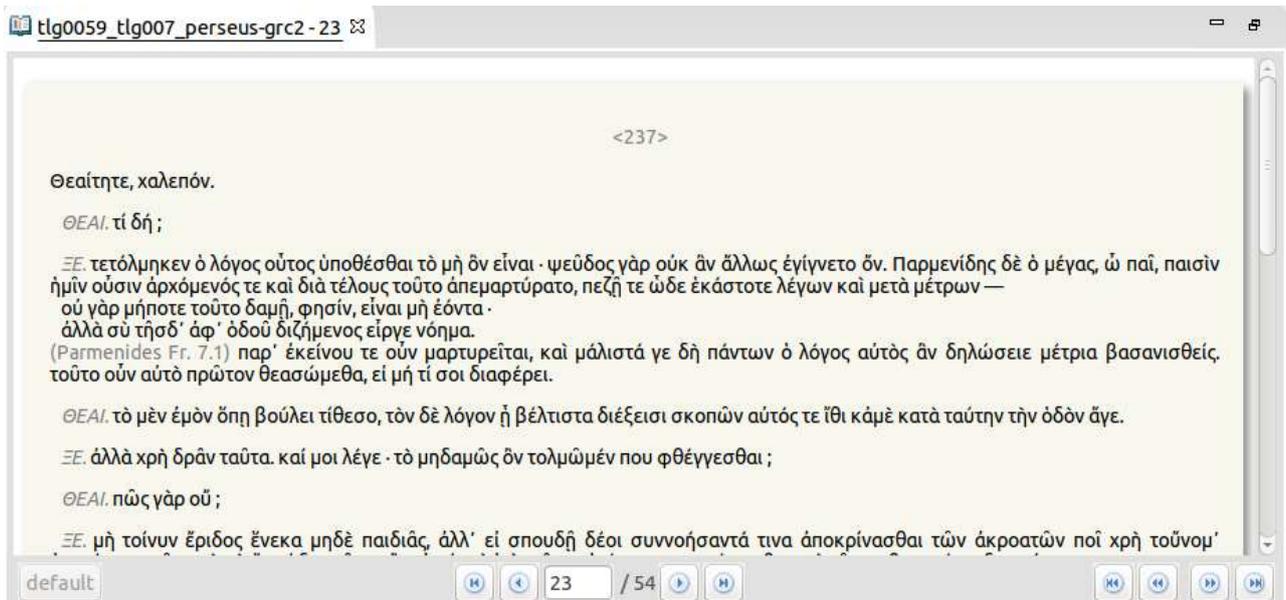


Figure 5. TXM edition view of *Sophist*, Estienne page 237a-e.

We should just stop a minute on that: it means that with this edition of Perseus corpus in TXM, one can search the word Γοργίας and differentiate between the instances of this word in the title (obviously one instance), in speaker names (97 instances in 1 dialogue),¹⁶ and in the speech of any protagonist (85 instances in 7 dialogues);¹⁷ one could also study which words are typical of Socrates' or Gorgias' speech. This functionality opens up a lot of possibilities, namely, to study the particular style of each protagonist of the platonic dialogue.

Bibliographic references encoded with <bibl> element in Perseus texts are declared and processed as a kind of note element for TXM XML-XTZ import: as such, they get a relevant display (see Figure 5) and are not mixed with the Greek content.

Every annotation made in Perseus texts, for instance, named entities (<name>, <persName>, <placeName>) or quotation tags (<q>, <quote>), are automatically available in TXM for search and analysis. An associated rendering could be added through the CSS stylesheet parameter if desired.

In addition, in order to test new possibilities offered by morphosyntactically tagged texts, we prepared a second corpus using the unique text of Plato available in Perseus TreeBank AGDT2:¹⁸

¹⁶ Example of query to find every beginning of Gorgias' speech turns:
<said>[_said_who="#Γοργίας" & _said_rend!="merge"]

¹⁷ Example of query to find mentions of Gorgias' name in the text: "Γοργί.*". The 7 dialogues with "Gorgias" occurrences are *Apology*, *Gorgias*, *Hippias Major*, *Meno*, *Phaedrus*, *Philebus*, and *Symposium*. Since we had no morphosyntactically tagged version of our texts, we use wild-cards to get a rough lemmatization. This query matches Γοργία, Γοργίας, Γοργίαν, Γοργίου, Γοργία.

¹⁸ Ancient Greek Dependency Treebank AGDT2: https://perseusdl.github.io/treebank_data.

Euthyphro. We imported this text in TXM via the XML/w loader with the Perseus TreeBank stylesheet¹⁹ as parameter. This corpus will be used here for one example in section 3.4.

3. A Typology of Textometric Analyses

For our example, we focus on the text of *Gorgias*, in comparison with Plato's other dialogues. So we define two structures inside the corpus: first, a *sub-corpus* containing only the *Gorgias* text, so that any textometric analysis (i.e. a KWIC concordance, a frequency list, etc.) can be processed on this text only. The second structure that we define is a *partition* of the corpus into texts, that is, we divide the corpus into parts (texts) so that we can compute contrastive analyses in order to compare texts to one another.

We have chosen to organize the typology according to the user's needs, that is, which kind of queries one would like to ask the corpus. These query types don't exactly match textometric functions, since one computation may be of use in several ways, and one type of query may get an answer through the complementary results of several functions.

3.1 Checking for occurrences and evaluating frequencies

To look at word frequencies in *Gorgias*, we can either compute raw frequencies, or compare the *Gorgias* subcorpus to the whole corpus of Plato using a statistic test, computing the specificities' score of each word (more precisely, each form of the word, since our corpus is not lemmatized) (see technical appendix 5.1).

Thus, for *Gorgias*' vocabulary, we can sort the results either by absolute frequency, or by specificities' score. In Figure 6, the left panel shows the most frequent words used in *Gorgias*, which here is useless since the most frequent words of *Gorgias* are the most frequent words in the corpus of Greek literature: they are “stop-words” like *καί*, *δέ*, and so on...²⁰ If we apply a stop-word filter,²¹ we get the most frequent content words of the text, giving account for the main vocabulary used in *Gorgias* (see middle panel).

¹⁹ `txm-filter-perseustreebank-xmlw.xsl`, available in TXM repository: <https://sourceforge.net/projects/txm/files/library/xsl/>.

²⁰ To see a list of the ancient Greek stop-words, cf. https://wiki.digitalclassicist.org/Stopwords_for_Greek_and_Latin.

²¹ There is no direct command in TXM to do this. If the corpus is tagged with part-of-speech, then this information makes it easy to select and filter out grammatical words (see example in Figure 18). Here, we downloaded Berra's stop-word list (<https://github.com/aurelberra/stopwords>); we used the `CreateCQPList` TXM macro (<https://groupes.renater.fr/wiki/txm-users/public/macros>) to read the stop-word file and build a variable containing all the stop-words (about 6,000 items); then we called this variable (“`grc_stopwords`”) in the query to

The figure displays three panels of word frequency analysis for the text *Gorgias*. Each panel shows a table with columns for 'word' and 'Frequency'.

Left Panel: Gorgias: [word!="\p{P}+"]

word	Frequency
καὶ	1473
τὸ	368
δὲ	354
ἢ	350
ὦ	281
τῶν	256
ἀν	244
γε	241
μὲν	236
εἶναι	226
ὅτι	214
γάρ	209
ὡς	189
τῆν	184
ἐν	183
τὰ	181
οὐ	180
τε	178
τοῦ	178
ὁ	176
μὴ	173
τὸν	169
οὐκ	168
οὖν	168
εἰ	167
περὶ	157
ἀλλ'	146
ἢ	141
τοῦς	137

Middle Panel: Gorgias: [word!=\$src_stopwords]

word	Frequency
Σώκρατες	104
δοκεῖ	92
Καλλίκλεις	62
δεῖ	50
ἀδικεῖν	48
οἶμαι	47
ἀγαθὸν	42
Πῶλε	39
ἀνάγκη	38
δικήν	36
ἀδικεῖσθαι	34
λόγον	33
Γοργία	32
ἀληθῆ	31
ἀνθρώπων	31
ὀρθῶς	31
δικαιον	30
πόλει	28
νυνδῆ	26
λόγος	25
οἶει	25
φαίνεται	25
λόγους	24
μέγα	24
βούλει	23
δηλον	23
σῶμα	23
δικαίως	22
Γοργίας	21
καλῶς	21

Right Panel: Gorgias: word

Units	Frequency T 613406	Gorgias t=30870	score
Καλλίκλεις	62	62	80.5
Πῶλε	39	39	50.6
Γοργία	32	32	41.5
ἀδικεῖσθαι	46	34	33.8
ἀδικεῖν	116	48	30.8
σὺ	797	124	27.8
ῥητορική	24	20	22.0
Γοργίας	28	21	21.3
ἐγὼ	944	124	21.2
ῥητορικής	23	19	20.8
αἴσιον	24	18	18.4
ἐστίν	976	120	18.2
σοι	953	116	17.3
ὦ	3262	281	17.2
;	9340	656	16.6
πῶλος	12	12	15.6
κάκιον	23	16	15.5
ῥητορικήν	25	16	14.7
φημί	43	20	14.5
ῥητορικήν	11	11	14.3
δικήν	169	36	12.7
ἐστίν	413	60	12.6
λέγω	458	64	12.6
ῥητορική	13	11	12.4
ἀδίκως	49	19	12.0
ῥητορική	11	10	12.0
ἔγωγε	462	63	11.9
φημι	56	20	11.8
ρήτωρ	18	12	11.4
ἦ	4776	350	11.4
διδόναι	64	20	10.6
Χαιρεφῶν	13	10	10.6
ἔλεγον	142	29	10.0
βελτίους	71	20	9.7
ἀγαθὸν	279	42	9.5
πιθανώτερος	7	7	9.1
πειθαίους	22	11	8.7

Figure 6. For *Gorgias*, the list of most frequent word forms²² (left panel), the same list without stop-words (middle panel), and the statistically most specific word forms (right panel).

The right panel may be more interesting, as it shows the words statistically over-used in *Gorgias* by comparison with the 28 other texts from Plato. The fact that the first three specific forms are Καλλίκλεις, Πῶλε, and Γοργία, that is, names of the protagonists, was expected, since those three protagonists do not appear elsewhere in Plato's work, and then are naturally specific to Plato's *Gorgias*.

exclude words matching any item of the stop- list. As TXM default tokenization preserves elided word forms (like ἀλλ', δ', etc.), for this experiment we just added to Berra's list the elided stop-words we found among the high-frequency words of our corpus. (We might also have changed the tokenization rules.)

²² TXM indexes punctuations too, so here the query filters punctuation signs out.

More significant is the fact that they are vocative forms, which are highly frequent in the dialog form: *Gorgias* is a dialog that is a real exchange (like *Laches*, *Euthydemus*, *Euthyphro*, *Hippias Major*, for instance), and not the very particular form of platonic dialog where the protagonists do not really interact intensively with each other, which are closer to a series of monologues (like the *Laws*, the *Parmenides*, *Republic*, and the *Timaeus*). This “dialogic” aspect of *Gorgias* became evident at the first sight of the specificities of the dialog: not only the over-use of vocative Σώκρατες, but also such dialog markers as σύ, ἐγώ, ὦ... Those first results show by linguistic means the refutative nature of the *Gorgias*, which is an attempt to refute sophistic positions and to expose the immoral implications in the hedonist position defended by Calicles. For Socrates, this requires repeatedly questioning the interlocutor (using the vocative), summarizing his argument (by using σὺ + verbs of saying φής, ὠμολόγεις, ἔλεγες, λέγεις...: “you says, you agree, you said...”), and opposing it with his own position (ἐγὼ λέγω, ἔλεγον...) (see also Figure 16 in section 3.4). Admittedly, *Gorgias* is not the only dialog where we can find those markers; they are common to what the scholarship calls the “Socratic dialogs,” which entails an ἔλεγχος, a refutation. But regarding those features, the textometric approach shows that *Gorgias* is representative, if not the dialog *par excellence* (Figure 7). It reveals that this dialog is paradigmatic, even something of an achievement of this method, since the dialog suddenly ends with a monolog and the acknowledgment that it is impossible to refute those who do not agree with the basic requirements of a rational discussion; that is why it ends with a myth. This analysis confirms the hypothesis that *Gorgias* is a transitional dialog, where Plato shows the limits of the Socratic method and the necessity to endorse another philosophical method.²³

²³ On this hypothesis, see (Marchand and Ponchon 2016:17–18).

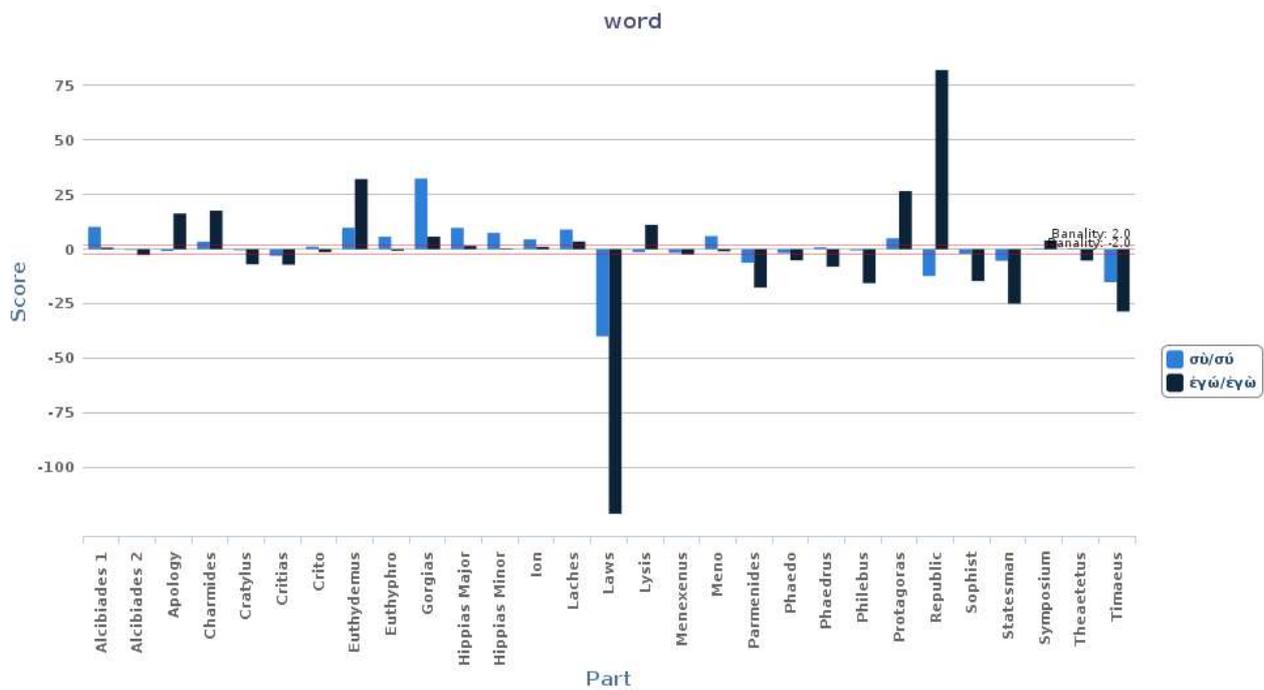


Figure 7. Specificities of ἐγώ and σύ in the texts of the Plato corpus.

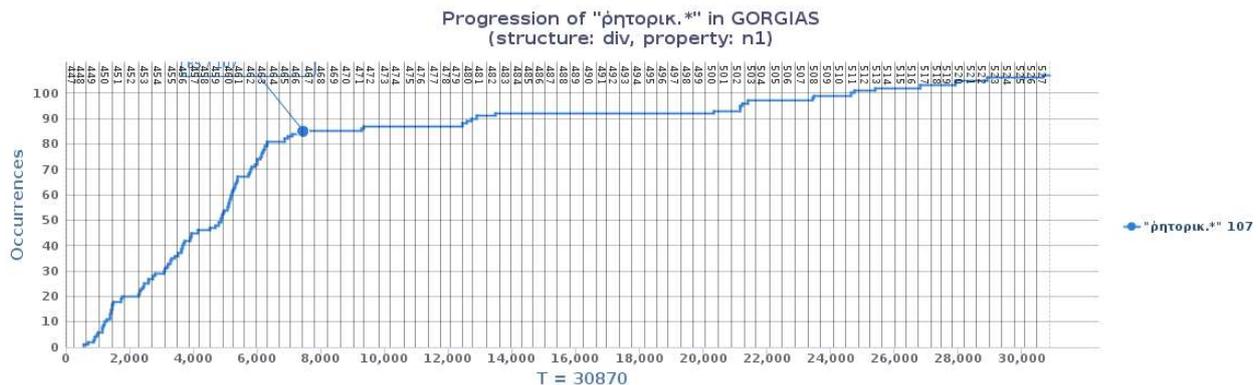
3.2 Visualizing the evolution of words or linguistic features

The TXM progression functionality is another way to analyze word distribution when data can be considered as a continuous evolution: this can be the case for a corpus when its texts are ordered according to time, or when focusing on one text, considering how it is organized from beginning to end. This second case can be applied to the *Gorgias*.

After the dialogic markers, the list of the specificities above shows also what are roughly the main themes of the dialog. The subtitle (which is probably not from the hand of Plato) of *Gorgias* is “on the rhetoric” (περὶ ῥητορικῆς), and it is well known that the main topic is a critique of the practices of the ῥήτορες, the men who deliver rhetorical speeches and teach rhetoric. Despite this very fact, the word ῥητορικῆ (more precisely, all the forms of the adjective ῥητορικός)²⁴ appear mostly in the first quarter of the dialog (there are 107 instances in *Gorgias*, and 85 appear from 448d to 466a, that is, in the discussion with Gorgias himself) (Figure 8).

²⁴ Since the corpus is not lemmatized, we submitted the query "ῥητορικ.*" (Figures 8 to 12 have been computed with this query).

GORGIAS: ["ῥητορικ.*"]



GORGIAS: "ῥητορικ.*"

Query: "ῥητορικ.*" Keyword: word Edit Search

sort keys: #1 None #2 None #3 None #4 None

1 - 100 / 107

ref	Left context	Keyword	Right context
Gorgias, s. 0463e	εις πολιτικῆς μορίου εἰδῶλον εἶναι τὴν	ῥητορικὴν	· ἀλλ' ἐγὼ τ
Gorgias, s. 0463e	ισομαι φράσαι ὃ γέ μοι φαίνεται εἶναι ἡ	ῥητορικὴ	· εἰ δὲ μὴ τ
Gorgias, s. 0465c	ὅτι ὁ ὄψοποικὴ πρὸς ἰατρικὴν, τοῦτο	ῥητορικὴ	πρὸς δικαιο
Gorgias, s. 0465d	ὄψοποικῶν. ὁ μὲν οὖν ἐγὼ φημι τὴν	ῥητορικὴν	εἶναι, ἀκήκ
Gorgias, s. 0466a	οὖν φῆς; κολακεία δοκεῖ σοι εἶναι ἡ	ῥητορικὴ	; κολακεία
Gorgias, s. 0467a	ῥήτορας νοῦν ἔχοντας καὶ τέχνην τὴν	ῥητορικὴν	ἀλλὰ μὴ κα
Gorgias, s. 0471d	σε ἐπήνεσα ὅτι μοι δοκεῖς εὐ πρὸς τὴν	ῥητορικὴν	παιδείουσ
Gorgias, s. 0471e	σοὶ ὡς ἐγὼ λέγω. ὦ μακάριε,	ῥητορικῶς	γάρ με ἐπυ
Gorgias, s. 0480a	Πῶλε, τίς ἡ μεγάλη χρεῖα ἐστὶν τῆς	ῥητορικῆς	; δεῖ μὲν γὰ

tlg0059_tlg023_perseus-grc2 - 22

<467>

τέχνην τὴν ῥητορικὴν ἀλλὰ μὴ κολακείαν, ἐμὲ ἐξελέγχας; εἰ δὲ με ἕσσεις ἀνέλεγκτον, οἱ ῥήτορες οἱ ποιούντες ἐν ταῖς πόλεσιν ἃ δοκεῖ αὐτοῖς καὶ οἱ τύραννοι οὐδὲν ἀγαθὸν τοῦτο κεκτήσονται, ἡ δὲ δύναμις ἐστίν, ὡς σὺ φῆς, ἀγαθόν, τὸ δὲ ποιεῖν ἀνευ νοῦ ἃ δοκεῖ καὶ σὺ ὁμολογεῖς κακὸν εἶναι· ἢ οὐ;

ΠΩΛ.

ἔγωγε.

ΣΩ.

πῶς ἂν οὖν οἱ ῥήτορες μέγα δύναντο ἢ οἱ τύραννοι ἐν ταῖς πόλεσιν, ἐὰν μὴ Σωκράτης ἐξελεγχθῆ ὑπὸ Πάλου ὅτι ποιοῦσιν ἃ βούλονται;

ΠΩΛ.

οὐτὸς ἀνὴρ—

ΣΩ.

default 22 / 82

Figure 8. Cumulative evolution graph for ῥητορικὴ in *Gorgias*²⁵ (upper frame), and hyperlinked concordance and text view on a selected occurrence (lower frames).

This is due to the fact that the Socratic criticism of rhetoric is a moral criticism of the way of life involved in the practice of the rhetoric. Socrates shows that the practice of rhetoric, as the faculty to persuade a crowd in order to take power, involves an immoral principle, namely, that it is better to do wrong than to suffer it. On the contrary, the philosophy of Socrates lies in the belief that it is better to suffer wrong than to do it (see *Gorgias* 469c, the question from Polos: εἰ δ' ἀναγκαῖον εἴη ἀδικεῖν ἢ ἀδικεῖσθαι, ἐλοίμην ἂν μᾶλλον ἀδικεῖσθαι ἢ ἀδικεῖν; "if it were necessary either to do wrong or to suffer it, I should choose to suffer rather than do it," translation Lamb from Perseus). For that reason, the forms ἀδικεῖσθαι and ἀδικεῖν are among the most specific forms of the *Gorgias* (they appear just after the three vocatives already

²⁵ X-axis represents the *Gorgias* text from its first word to its last word; Y-axis represents the number of occurrences of words beginning by ῥητορικ- since the beginning of the text. Thus, the more the curve rises, the more occurrences there are in the current passage. A flat curve means no occurrences of the word in the passage.

mentioned Καλλίκλεις, Πῶλε, and Γοργία), although they are in general quite common in Plato's work.²⁶

A glance at the table of frequencies and graph of specificities of the adjective ῥητορικός (all its cognates) in every text of our corpus (Figure 9) shows, however, that ῥητορικός is typical of the *Gorgias*, for the very reason that it is not very common in Plato's work: hence, even if Plato does not continually use the word in the dialog, it is nonetheless characteristic of the *Gorgias* (107 instances, specificity of +108), and to a lesser extent the *Phaedrus* (21 instances, specificity of +8) (these texts are the only two where the specificity score is greater than +3).²⁷

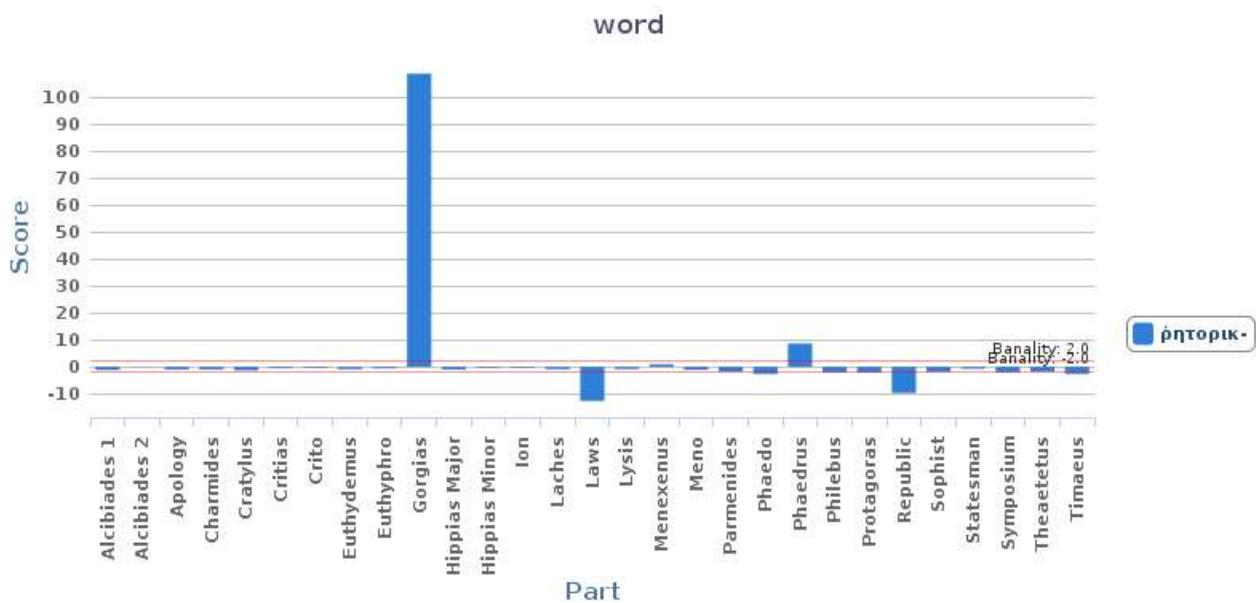


Figure 9. Specificities of ῥητορική in the texts of the Plato corpus.

E. Schiappa has polemically stated that it could be possible that Plato coined the word ῥητορική.²⁸ This is not the place to take a position on the external evidence discussed by Schiappa, who claims that all the texts that employ the word ῥητορική in the fourth century can be considered contemporary or later than Plato.²⁹ But the tool here allows us to confirm the fact that in Plato “the various forms of ῥητορική are curiously distributed.” (Schiappa 1990:463) and

²⁶ 46 instances in Plato of ἀδικεῖσθαι (34 in *Gorgias*); 116 instances of ἀδικεῖν (48 in *Gorgias*). The query "ἀδικ.*" on whole of Plato's work shows, with no surprise, a massive presence in *Republic* (164 instances), *Gorgias* (163 instances), and *Laws* (103 instances); the presence in *Republic* and *Laws* has to be relativized by the size of the dialogs (respectively 104,572 and 116,925 words, whereas for *Gorgias* there are 30,870 words).

²⁷ See the other instances, *Statesman* 304e, 304d; *Republic* 548e; *Menexenus* 235e, 236a (bis), *Euthydemus* 307a, *Cratylus* 425a, *Alcibiades 2* 145e, *Thaetetus* 177b.

²⁸ Schiappa 1990. See also the criticisms of O'Sullivan 1993 and the revised Schiappa 1999.

²⁹ By the facts, the dates of those texts are difficult to determine, as is the case for Alcidamas or the *Rhetoric to Alexander*.

the graph shows that this curious distribution relies on the specificity of *Gorgias*' vocabulary (Figure 9).³⁰

3.3 Refining a word's meaning with systematic contextual use

The KWIC concordance view is a core functionality in the textometric approach, as it shows very precisely and efficiently how a word is used in the text. The context size can be adjusted if needed, but for a view of a larger selection of text, a double-click on a concordance line opens a hypertext link to the corresponding text page with the search word highlighted. The two views are dynamically aligned: selecting a line in the concordance shows the corresponding page in the text edition (Figure 10).

ref	Left context	Keyword	Right context
Phaedrus, s. 0271d	γχάνει ψυχαγωγία ούσα, τὸν μέλλοντα	ρητορικὸν	ἔσσεσθαι ἀνάγκη εἶδέναι ψυ
Phaedrus, s. 0272d	ἰδόντων ἢ τροφῆ, τὸν μέλλοντα ἰκανῶς	ρητορικὸν	ἔσσεσθαι. τὸ παράπαν γὰρ οἱ
Alcibiades 2, s. 0145e	τὸ ἀποκτείνουσι, πρὸς δὲ καὶ ἀνδρῶν	ρητορικῶν	πολιτικῶν φύσημα φυσῶνται
Euthydemus, s. 0307a	ὃν δοκεῖ σοὶ εἶναι, καὶ χρηματιστικὴ καὶ	ρητορικὴ	καὶ στρατηγία; ἔμοιγε πάντι
Gorgias, s. 0448d	ς καὶ ἐξ ὧν εἰρήκεν ὅτι τὴν καλουμένην	ρητορικὴν	μᾶλλον μεμελέτηκεν ἢ διαλέ
Gorgias, s. 0449a	καλεῖν ὡς τίνος ἐπιστήμονα τέχνης. τῆς	ρητορικῆς	, ὡς Σώκρατες, ῥήτορα ἄρα:
Gorgias, s. 0449c	βραχυλογωτέρου ἀκούσαι. φέρε δὴ·	ρητορικῆς	γὰρ φῆς ἐπιστήμων τέχνης ε
Gorgias, s. 0449d	καὶ ποιῆσαι ἂν καὶ ἄλλον ῥήτορα· ἢ	ρητορικῆ	περὶ τῶν ὄντων τυγχάνει ε
Gorgias, s. 0449d	θι δὴ μοι ἀπόκριται οὕτως καὶ περὶ τῆς	ρητορικῆς	, περὶ τῶν ὄντων ἐστὶν ἐπι
Gorgias, s. 0449e	οὐκ ἄρα περὶ πάντας γε τοὺς λόγους ἢ	ρητορικῆ	ἐστίν. οὐ δὴτα. ἀλλὰ μὴν λέγ
Gorgias, s. 0450b	τί οὐδὲν ποτε τὰς ἄλλας τέχνας οὐ	ρητορικὰς	καλεῖς, οὐσας περὶ λόγους, ε
Gorgias, s. 0450b	αλεῖς, οὐσας περὶ λόγους, εἶπερ ταύτην	ρητορικὴν	καλεῖς, ἢ ἂν ἢ περὶ λόγους; ἔ
Gorgias, s. 0450b	εἰπεῖν πᾶσα ἐστὶν ἢ ἐπιστήμη, τῆς δὲ	ρητορικῆς	οὐδὲν ἐστὶν τοιοῦτον χειρο
Gorgias, s. 0450c	διὰ λόγων ἐστίν. διὰ ταῦτ' ἐγὼ τὴν	ρητορικὴν	τέχνην ἀξιώ εἶναι περὶ λόγου
Gorgias, s. 0450d	δοκεῖς λέγειν, περὶ ἧς οὐ φῆς τὴν	ρητορικὴν	εἶναι· ἢ οὐ; πάνυ μὲν οὐν κα
Gorgias, s. 0450e	ὧν τοιούτων τινὰ μοι δοκεῖς λέγειν τὴν	ρητορικὴν	ἀληθῆ λέγεις. ἀλλ' οὗτοι το
Gorgias, s. 0450e	οὔτως γε οὐδεμίαν οἰμαί σε βούλεσθαι	ρητορικὴν	καλεῖν, οὐκ ὅτι τῷ ῥήματι οἰ
Gorgias, s. 0450e	, ὅτι ἢ διὰ λόγου τὸ κύρος ἔχουσα	ρητορικῆ	ἐστίν, καὶ ὑπολάβοι ἂν τις, ε
Gorgias, s. 0450e	εἶν ἐν τοῖς λόγοις, τὴν ἀριθμητικὴν ἄρα	ρητορικὴν	, ὡς Γοργία, λέγεις; ἀλλ' οὐκ

Figure 10. Concordance of *ρητορικὴ* in the Plato corpus, and the hyperlinked view of the text page corresponding to one selected concordance line.

Furthermore, the table layout of the KWIC view, coupled with the possibility to sort right and left contexts, reveals the lexical patterns involving the searched word. For instance, sorting on the 'right context' of *ρητορικὴ* allows one to find the passages where the concept is defined: all the contexts where *ρητορικὴ* is followed by the verb *εἶμι* (to be) or *καλέω* (to call) (Figure 11).

³⁰ By the way, among the 30 first words with the highest score of specificity in *Gorgias*, there are 6 forms of ἢ *ρητορικὴ*.

Query: ῥητορικῆ Keyword: word Edit Search

sort keys: #1 Right cont #2 None #3 None #4 None Sort

1 - 100 / 138 >

ref, said_who	Left context	Keyword	Right context
Gorgias, s. 0460a, #Σωκράτης	ς, ὥσπερ ἄρτι εἶπες, ἀποκαλύψας τῆς	ῥητορικῆς	εἰπέ τις ποθ' ἡ δύναμις ἐστίν. ἀλλ'
Gorgias, s. 0462c, #Πῶλος	περγασίας. οὐκοῦν καλόν σοι δοκεῖ ἢ	ῥητορικῆ	εἶναι, χαρίζεσθαι οἷόν τε εἶναι ἀνέ
Gorgias, s. 0465d, #Σωκράτης	ὄψοποικῶν. ὁ μὲν οὖν ἐγὼ φημι τὴν	ῥητορικὴν	εἶναι, ἀκήκοας· ἀντίστροφον ὄψο
Phaedrus, s. 0269d, #Σωκράτης	ἄπερ τάλλα· εἰ μὲν σοι ὑπάρχει φύσει	ῥητορικῶ	εἶναι, ἔση ῥήτωρ ἐλλόγιμος, προσ
Gorgias, s. 0462b, #Πῶλος	τάληθ' εἰρήσθαι. ἀλλὰ τί σοι δοκεῖ ἢ	ῥητορικῆ	εἶναι; πράγμα ὁ φῆς σὺ ποιῆσαι τ
Gorgias, s. 0462c, #Πῶλος	ἐγώγε τινα. ἐμπειρία ἄρα σοι δοκεῖ ἢ	ῥητορικῆ	εἶναι; ἔμοιγε, εἰ μή τι σὺ ἄλλο λέγε
Gorgias, s. 0450d, #Σωκράτης	δοκεῖς λέγειν, περὶ ἧς οὐ φῆς τὴν	ῥητορικὴν	εἶναι· ἢ οὐ; πάνυ μὲν οὖν καλῶς ὑ
Gorgias, s. 0450e, #Σωκράτης	πτων γε οὐδεμίαν οἰμαί σε βούλεσθαι	ῥητορικὴν	καλεῖν, οὐκ ὅτι τῷ ῥήματι οὕτως ε
Gorgias, s. 0450b, #Σωκράτης	τί οὖν δὴ ποτε τὰς ἄλλας τέχνας οὐ	ῥητορικὰς	καλεῖς, οὐσας περὶ λόγους, εἶπερ
Gorgias, s. 0450b, #Σωκράτης	λεῖς, οὐσας περὶ λόγους, εἶπερ ταύτην	ῥητορικὴν	καλεῖς, ἢ ἂν ἢ περὶ λόγους; ὅτι, ὡ
Phaedrus, s. 0260c, #Σωκράτης	ποῖόν τιν' ἂν οἶει μετὰ ταῦτα τὴν	ῥητορικὴν	καρπὸν ἂν ἔσπειρε θερίζειν; οὐ γ
Gorgias, s. 0463d, #Σωκράτης	γμάθοις ἀποκριναμένου; ἔστιν γὰρ ἢ	ῥητορικῆ	κατὰ τὸν ἐμὸν λόγον πολιτικῆς με
Gorgias, s. 0453c, #Σωκράτης	ἵνα ποτε λέγεις τὴν πειθῶ τὴν ἀπὸ τῆς	ῥητορικῆς	καὶ περὶ τίνων αὐτὴν εἶναι. τοῦ ἐν
Euthydemus, s. 0307a, #Σωκράτης	δοκεῖ σοι εἶναι, καὶ χρηματιστικῆ καὶ	ῥητορικῆ	καὶ στρατηγία; ἔμοιγε πάντως δὴ
Gorgias, s. 0502d, #Σωκράτης	· τί δὲ ἢ πρὸς τὸν Ἀθηναίων δῆμον	ῥητορικῆ	καὶ τοὺς ἄλλους τοὺς ἐν ταῖς πόλε
Gorgias, s. 0461a, #Σωκράτης	ὄν ἀδύνατον εἶναι ἀδίκως χρῆσθαι τῇ	ῥητορικῆ	καὶ ἐθέλειν ἀδικεῖν. ταῦτα οὖν δη
Gorgias, s. 0461a, #Σωκράτης	ἰλίγον ὑπερτον ἔλεγες ὅτι ὁ ῥήτωρ τῇ	ῥητορικῆ	κὰν ἀδίκως κράτο, οὕτω θαυμά
Gorgias, s. 0450e, #Σωκράτης	τὴν ἀριθμητικὴν οὔτε τὴν γεωμετρικὴν	ῥητορικὴν	λέγειν. ὀρθῶς γὰρ οἶει, ὡ Σώκρα
Phaedrus, s. 0263b, #Σωκράτης	ῥώμεθα. οὐκοῦν τὸν μέλλοντα τέχνην	ῥητορικὴν	μετιέναι πρῶτον μὲν δεῖ ταῦτα ὀ

default 17 / 82

Figure 11. Concordance of ῥητορικῆ in the Plato corpus, sorted on 'right context'.

Lastly, to have an overview of the platonic approach of the word ῥητορικῆ, one can have a look at its cooccurrents³¹ in the *Gorgias*: this query shows the words especially used by Plato in proximity with ῥητορικῆ and its cognates (Figure 12). The cooccurrence association score is computed with the same statistic measure as the specificity explained above (the cooccurrence score is precisely the specificity of the cooccurrent word in the part formed by the whole set of contexts of ῥητορικῆ in comparison with its global frequency in the corpus). The common feature of the platonic definition of rhetoric appears within a range from -10 to +10 words:³² rhetoric is the name of the competence of the ῥήτωρ; Socrates defines it as “producer of persuasion” (e.g. 453a: πειθοῦς δημιουργός, 453d: πειθῶ ποιεῖν); its definition is given through an analogy as part (μόριον) of the art of flattery (κολακεία; e.g. 463b); it is used in the context of the court (δικαστηρίοις)... It is worth noting that those cooccurrents are also the main cooccurrents of the word ῥητορικῆ in all Plato's works.

³¹ What is called “cooccurrence” in textometry is often called “collocation” in corpus linguistics.

³² This context size is a common choice, but the software allows for the exploration of different settings parameters if needed.

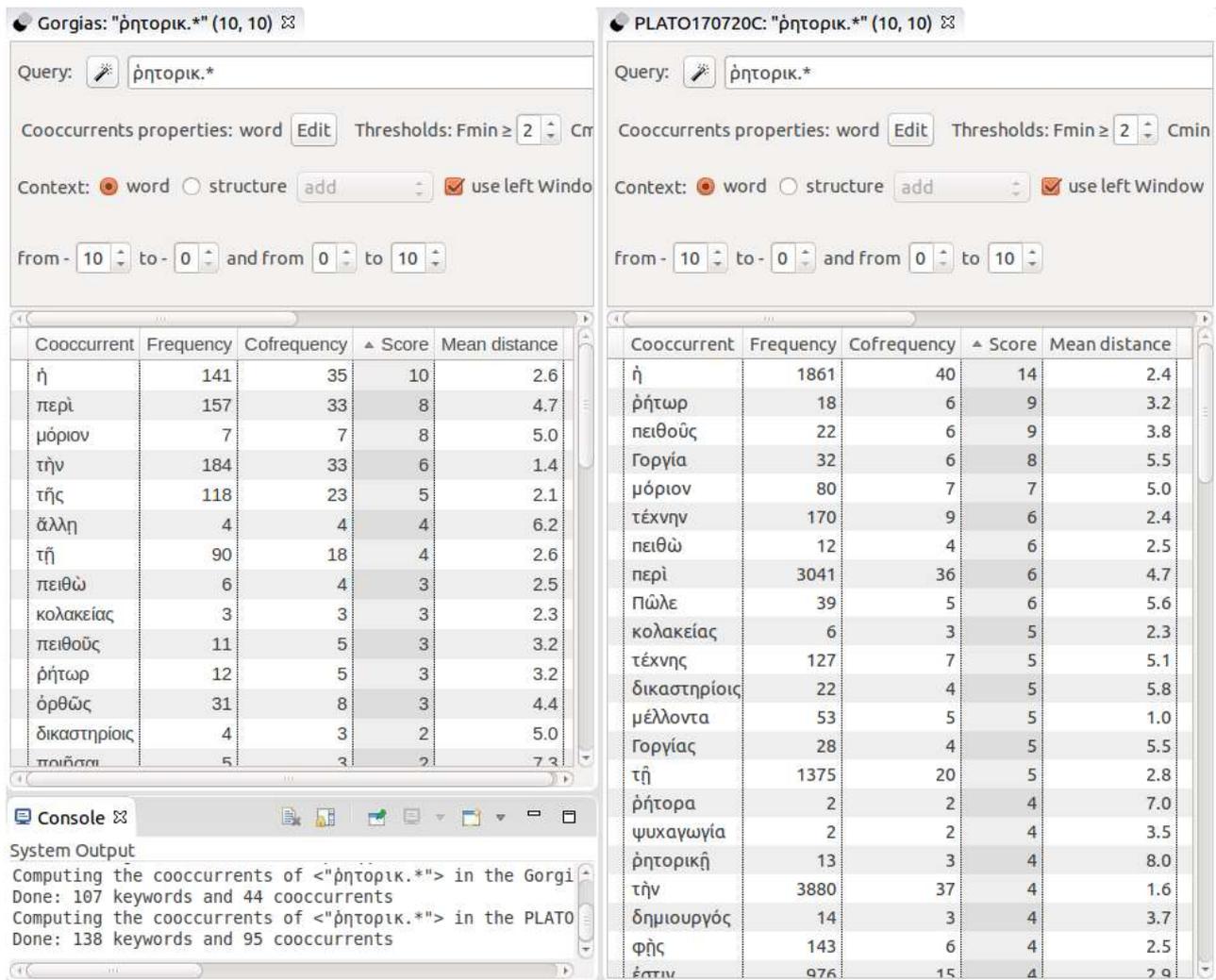


Figure 12. Cooccurrences of ῥητορικῆ in *Gorgias* (left panel) and in the whole Plato corpus (right panel).

In the same fashion, the query of the cooccurrences of ἀδικε.* (which matches both ἀδικεῖσθαι and ἀδικεῖν) (Figure 13) instantly points to the passages where Plato makes the comparison between suffering and doing wrong. It is no surprise that the cooccurrences are ἀδικεῖσθαι (because it is most of the time compared to ἀδικεῖν), ἀδικεῖν (for the same reason), τὸ (because the comparison is between the fact of doing or suffering wrong), αἴσχιον, and κακίον (because in the dialog with Polos, the question is mainly whether it is fouler or more evil to commit or to suffer wrongdoing).

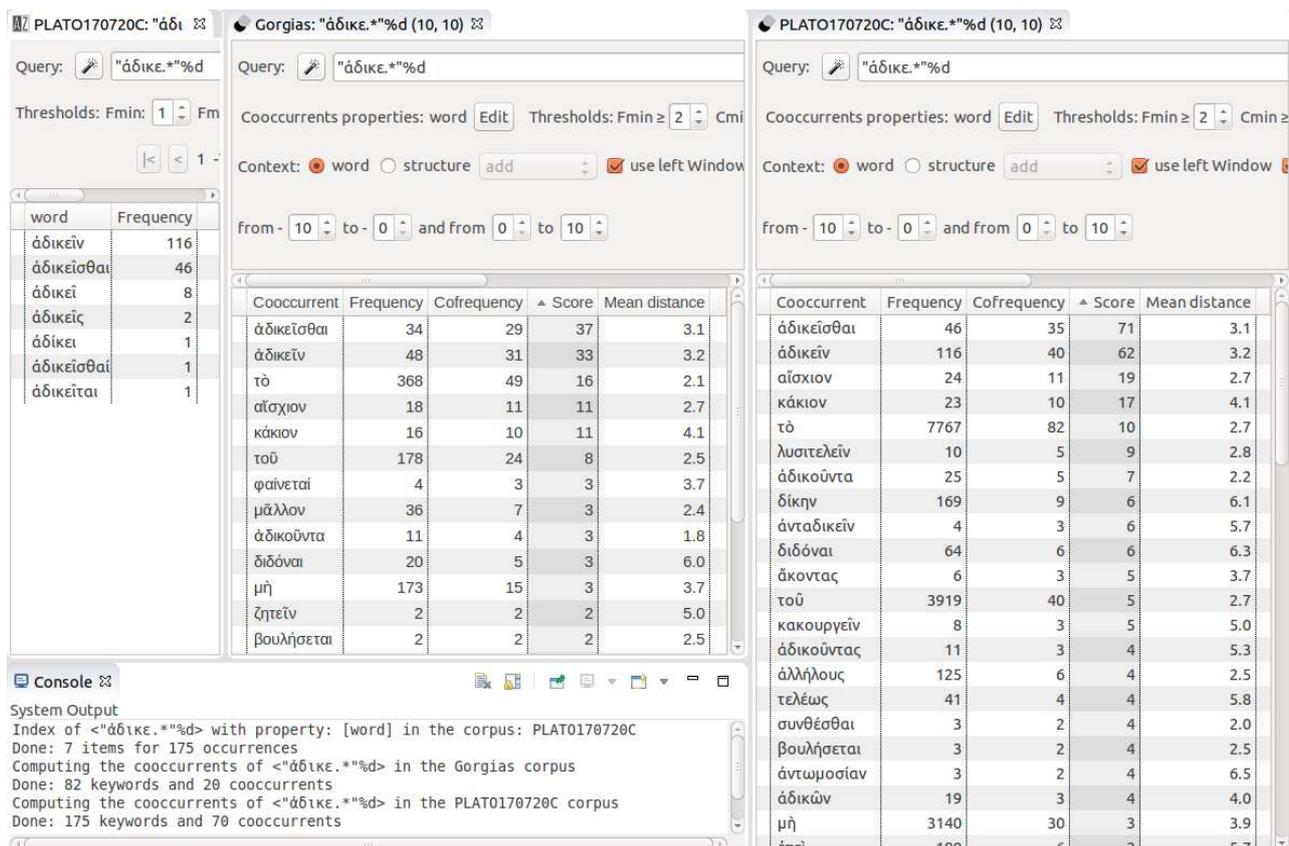


Figure 13. Cooccurrences of ἀδικεῖσθαι and ἀδικεῖν (cf. complete word forms list on the left) in *Gorgias* (central panel) and in the whole Plato corpus (right panel).

An interesting case arises when cooccurrences in the *Gorgias* happen to be quite different from the cooccurrences in the whole corpus, that is, instances in which a word gets a distinct meaning in the *Gorgias*. This is the case for νόμος (Figure 14): in the *Gorgias* ὁ νόμος is mainly used in the discussion with Callicles in the context of the opposition between what is κατὰ νόμον (“according to convention”) and κατὰ φύσιν (“according to nature”). In the other dialogues (and mainly in the *Laws*) the cooccurrences of νόμος emphasize the description of what will be the law, or conform to the law (ἔστω κατὰ νόμον), or the political action of the law (κείσθω νόμος).

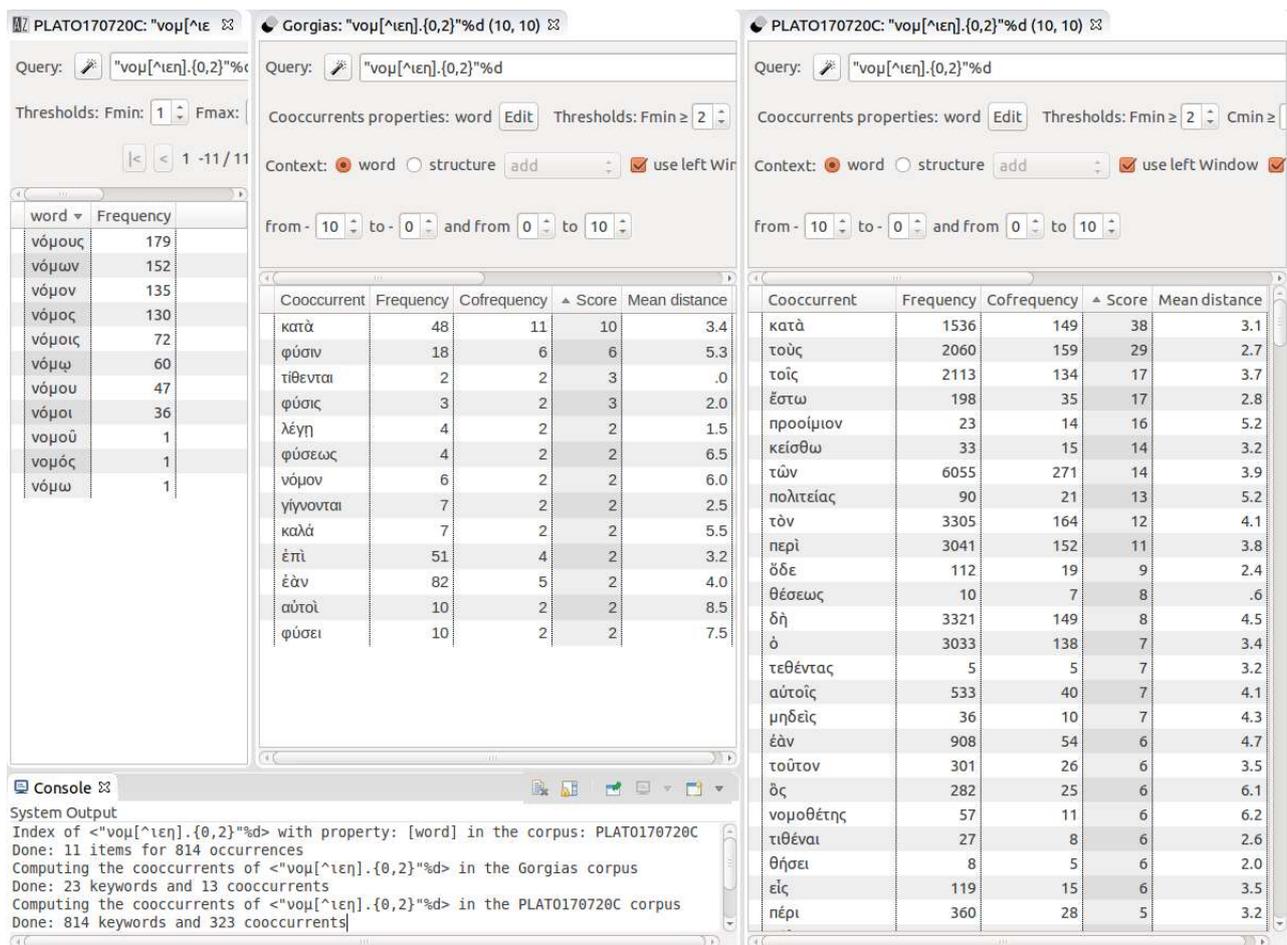


Figure 14. Cooccurrences of νόμος (cf. complete word forms list on the left) in *Gorgias* (central panel) and in the whole Plato corpus (right panel).

3.4 Computing corpus-based paradigmatic series

Another way of using the search engine is to build a list of words or patterns that share some constitutive or contextual feature. For instance, one could list all the lexical units including the same given morphological unit, or list all units occurring in some precise position (words ending verses, adjectives qualifying a given noun).

Schiappa emphasized “Plato was a prolific coiner of words ending with -ική denoting ‘art of’” (Schiappa 1990:464). We can have a look at our data to get the set of these kinds of words (Figure 15).

Query: Properties: word Edit Search

Thresholds: Fmin: 1 Fmax: 9999999 Vmax: 9999999 Page size: 100

1 -56 / 56 t120 , v56 , fmin1 , fmax17

word	Frequency
ίατρική	17
ρήτορική	11
γυμναστική	6
μμητική	5
μουσική	5
λογιστική	4
σκυτοτομική	4
ύπηρετική	4
πολεμική	3
πολιτική	3
σοφιστική	3
Κρητική	2
αύλητική	2
βασιλική	2
γραφική	2
διακριτική	2
μετρική	2
στρατηγική	2
τεκτονική	2
άριθμητική	2
έμπορική	2

Figure 15. Index for the query *.ική in the 29 texts of Plato, results for a minimum frequency of 2.

In section 3.1 above, we noted that in the *Gorgias*, Socrates frequently uses the pronoun σύ with verbs of saying, and ἐγώ with verbs expressing his own position. We can get a rough but systematic summary of verbs more frequently used with ἐγώ and σύ in the *Gorgias* by computing their cooccurrences in a narrow context window (5 words right and left) (Figure 16).

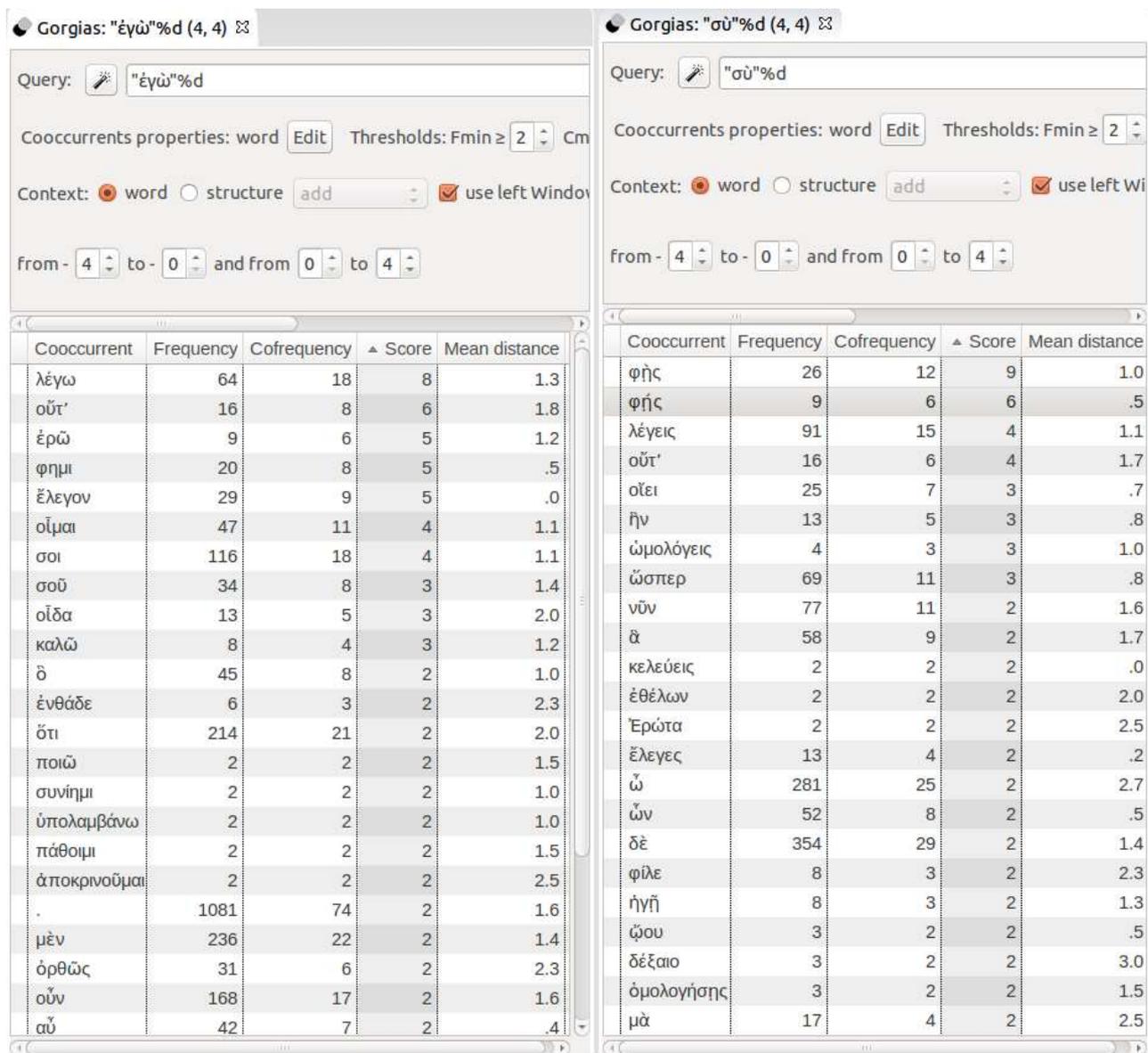


Figure 16. Cooccurrences of ἐγώ and σύ in Plato's *Gorgias*.

Another kind of paradigmatic set can be computed on the basis of textual structures. For instance, we can have a look at every one-word answer in the corpus, and thus explore the different ways in Plato's dialogues to express agreement or disagreement (Figure 17).

PLATO170720C: (<said>[]</said> | <sp>[]</sp>):word

Query: (<said>[]</said> | <sp>[]</sp>)| Properties: word Edit Search

Thresholds: Fmin: 1 Fmax: 9999999 Vmax: 9999999 Page size: 100

1 -100 / 172 t1377 , v172 , fmin1 , fmax430

word	Frequency
ναί .	430
ἔγωγε .	87
πῶς ;	72
φαίνεται .	68
οὕτως .	50
ὀρθῶς .	47
ἀνάγκη .	43
ἔοικεν .	41
οὐδαμῶς .	32
ἀληθῆ .	30
φημί .	27
ἀληθέστατα .	27
ἔμοιγε .	26
ἴσως .	21
πῆ ;	20
ὀρθότατα .	17
κινδυνεύει .	16
καλῶς .	15
οὐδέν .	14
δῆλον .	12
ἔστιν .	12

Figure 17. Index of one-word answers.

But the Plato TXM corpus we are using here does not reveal all the possibilities given by this kind of query. A more precise investigation could be done on a tagged corpus. Let us have a quick look to two other examples taken outside of *Gorgias*.

Our first example deals with morphosyntactic annotation. The *Euthyphro* corpus was built from the TreeBank AGDT2 and has morphosyntactic tags, so we can list the verbs occurring in the same sentences as πατήρ (Figure 18).

The screenshot shows the TXM software interface with the following components:

- Search Bar (Upper Left):** Query: [postag="n.*"];lemma. Threshholds: Fmin: 1, Fmax: 9999999. Results table:

lemma	Frequency
θεός	58
Σωκράτης	41
Ευθύφρων	38
δίκη	18
πατήρ	16
άνθρωπος	12
λόγος	11
Ζεύς	10
θεοσσία	10
- Search Bar (Upper Central):** Query: [postag="v.*"]; Threshholds: Fmin: 1, Fmax: 9999. Results table:

lemma	Frequency
ειμί	9
ποιέω	7
ἐπέειμι	5
λέγω	4
τυγχάνω	4
διδάσκω	3
συνδέω	3
ἔχω	3
- Page View (Upper Right):** Greek text snippet:

τὸ δσιον καὶ τί τὸ ἀνόσιον ; [0]
λέγω τοίνυν ὅτι τὸ μὲν δσιόν ἐστιν ὄπερ ἐγὼ νῦν ποιῶ, τῷ
ἀδικοῦντι ἢ περὶ φόνου ἢ περὶ λερῶν κλοπᾶς ἢ τι ἄλλο τῶν
τοιούτων ἐξαμαρτάνοντι ἐπεξιέναι, ἔαντε πατήρ ὦν τυγχάνη
ἔαντε μήτηρ ἔαντε ἄλλος ὅστισούν, τὸ δὲ μὴ ἐπεξιέναι ἀνόσιον .
[0]
ἐπεὶ ὦ Σώκρᾳτες, θέασαι ὡς μέγα σοι ἐρῶ τεκμήριον τοῦ νόμου
ὅτι οὕτως ἔχει- ὁ καὶ ἄλλοις ἤδη εἶπον, ὅτι ταῦτα ὀρθῶς ἂν εἴη
οὕτω γιγνόμενα- μὴ ἐπιτρέπειν τῷ ἀσεβοῦντι μὴδ ἂν ὅστισούν
τυγχάνη ὦν.
αὐτοὶ γὰρ οἱ ἄνθρωποι τυγχάνουσι νομίζοντες τὸν Δία τῶν θεῶν
ἄριστον καὶ δικαιοτάτον, καὶ τοῦτον ὁμολογοῦσι τὸν αὐτοῦ
πατέρα δῆσαι ὅτι τοὺς υἱεὶς κατέπινεν οὐκ ἐν δίκῃ, κάκεινόν γε
αὐτὸν αὐτοῦ πατέρα ἐκτεμνείν δι' ἕτερα τοιαῦτα .
ἐμοὶ δὲ χαλεπαίνουσι ὅτι τῷ πατρὶ ἐπέξερχομαι ἀδικοῦντι, καὶ
οὕτως αὐτοὶ αὐτοῖς τὰ ἐναντία λέγουσι περὶ τε τῶν θεῶν καὶ
περὶ ἐμοῦ.
ἀλλ' οὐκ ἐθέλωμαι τὰ ἐναντία λέγειν αὐτοῖς τῶν νοσούντων ἀνθρώπων .
- Search Bar (Lower):** Query: [lemma="πατήρ"]; Keyword: word. Sort keys: #1 Left conte, #2 None, #3 None, #4 None. Results table:

text_id	Left context	Keyword	Right context
tlg0059.tlg001.perseus-grc1.tb	, κάκεινόν γε αὐτὸν αὐτοῦ	πατέρα	ἐκτεμνείν δι' ἕτερα τοιαῦτα . ἐμοὶ δὲ χαλεπαίνουσι ὅτι τῷ
tlg0059.tlg001.perseus-grc1.tb	διαφθεῖροντι ἐμέ τε καὶ τὸν αὐτοῦ	πατέρα	, ἐμὲ μὲν διδάσκοντι, ἐκείνον δὲ νουθετοῦντί τε καὶ
tlg0059.tlg001.perseus-grc1.tb	, καὶ τοῦτον ὁμολογοῦσι τὸν αὐτοῦ	πατέρα	δῆσαι ὅτι τοὺς υἱεὶς κατέπινεν οὐκ ἐν δίκῃ, κάκεινόν
tlg0059.tlg001.perseus-grc1.tb	ἡμετέρων ἀποσφάττει αὐτόν. ὁ οὖν	πατήρ	συνδήσας τοὺς πόδας καὶ τὰς χεῖρας αὐτοῦ, καταβαλὼν εἰς
tlg0059.tlg001.perseus-grc1.tb	ἐπεχείρησας ὑπὲρ ἀνδρὸς θητὸς ἀνδρα πρῆσβύτην	πατέρα	δικωκάθειν φόνου, ἀλλὰ καὶ τοὺς θεοῦς ἂν ἐδεισας παρακινδυνεύειν

Figure 18. Verbs used with πατήρ in *Euthyphro*. (i) (upper left) search for most frequent nouns in the text (lemmas); (ii) (lower) contexts of πατήρ in a concordance view; (iii) (upper right) widening context in a page view; (iv) (upper central) list of verbs (lemmas) in sentences containing πατήρ, sorted by descending frequency.

Obviously, all encoded information can be used to define what will be returned. Our second example uses semantic information. In our Plato corpus, one text (the *Phaedo*) has been tagged for named entities; we can use this annotation to see precisely which places or people are mentioned in *Phaedo* (Figure 19).

The figure displays three screenshots of the Perseus search interface, each showing the results of a query for named entities in the *Phaedo* text. The queries are: 1) `<name>[]+</name>:w`, 2) `<persname>[]+</persname>`, and 3) `<placename>[]+</placename>`. Each screenshot shows a table with 'word' and 'Frequency' columns.

word	Frequency
Ἀθηναῖοι	3
Ἀθηναίων	2
Ἀχερουσιάδα	2
Ἀύριον	1
Βαβαί	1
Βοιωτοῦς	1
Θηβαϊκῆς	1
Θηβαῖε	1
Θηβαῖος	1
Παιανιεὺς	1
Στύγιον	1
Φλειασίων	1
Ἀθηναίσις	1
Ἄργεῖοι	1
Ἀχερουσιάδι	1
Ἀχερουσιάδος	1
Ἡρακλείων	1
Ἰόλεων	1

word	Frequency
Κέβης	82
Σώκρατες	65
Σιμμίας	38
Σωκράτης	35
Σιμμία	34
Κρίτων	16
Ἄϊδου	15
Δία	12
Φαίδων	12
Σωκράτους	8
Σιμμίαν	7
Ἐχέκρατες	7
Σωκράτη	6
Τάρταρον	6
Σιμμίου	5
Κρίτωνα	4
Δί '	3
Εὐήνω	3
Εὐήνος	3

word	Frequency
Δῆλον	2
Πυριφλεγέθοντα	2
Πυριφλεγέθοντι	2
Αἰγίνη	1
Αἰγύπτω	1
Δήλου	1
Δηλόν	1
Εὐρίπω	1
Κρήτην	1
Κωκυτόν	1
Μέγαρα	1
Μεγαρόθεν	1
Σικελία	1
Στύγα	1
Ἀθήναζε	1
Ἄτλαντα	1
Ἑλλάς	1
Ὀδυσσεΐα	1

Figure 19. Index of named entities encoded in the Perseus edition of *Phaedo*.

3.5 Local contrastive analysis of a corpus: Identifying what is typical in a part

Studying the vocabulary of *Gorgias*, we have already shown that the specificity measure can be used according to two points of view: for a word (ex. ῥητορική, ἐγή, σύ) we can draw its usage profile among texts; and for a text, we can point out the words that are statistically overused in this text in comparison to the rest of the corpus. This way, it is possible to study and compare every division of the corpus. For example, if we have a hypothesis that a corpus is divided into four main parts (be it text types, time spans, streams, etc.), we have a tool to automatically extract words that could help grounding and exploring this hypothesis.

As a complement to statistical processing, another way to look at the particularities of the *Gorgias*' vocabulary is by simply listing the words that appear only in *Gorgias*³³ (Figure 20).

³³ In TXM, the table used is the one produced by computing the specificities on the *Gorgias* subcorpus

Units	Frequency T 613406	Gorgias t=30870	score	PLATO170720C \ Gorgias
Καλλίκλεις	62	62	80.5	0
Πῶλε	39	39	50.6	0
Γοργία	32	32	41.5	0
πῶλος	12	12	15.6	0
ῥητορικὴν	11	11	14.3	0
πιθανώτερος	7	7	9.1	0
Καλλικλῆς	6	6	7.8	0
Πῶλον	6	6	7.8	0
ὄψοποικὴ	5	5	6.5	0
κόρρης	4	4	5.2	0
ἐλέγχου	4	4	5.2	0
ἐπιψηφίζειν	4	4	5.2	0
ἐροῦ	4	4	5.2	0
δημηγορία	3	3	3.9	0
διδούς	3	3	3.9	0
ἠύρισκομεν	3	3	3.9	0
κάεσθαι	3	3	3.9	0
προσφέρει	3	3	3.9	0
ἀπαλλαττόμενος	3	3	3.9	0
ἀποκτενεῖ	3	3	3.9	0
ἀφαιρήσεται	3	3	3.9	0
Ἀρχέλαος	3	3	3.9	0
ἑάσεις	3	3	3.9	0
ἔλεγές	3	3	3.9	0
ἔλεγε	3	3	3.9	0
ὄψοποιία	3	3	3.9	0
ὁμολογήσης	3	3	3.9	0

Figure 20. Word forms used only in Plato's *Gorgias*, with a minimum frequency of 3 (null frequency in all other texts of the Plato corpus).

For instance, putting aside proper names, the word-form ἐπιψηφίζειν (to put a question to the vote) has 4 instances in Plato, all of them in *Gorgias*—to which we should add one instance of ἐπιψηφίζων (476a). Even if the term is rather technical, *prima facie* this result might be surprising, given the fact that Plato is known as a critic of Athenian democracy and of the fact that the citizens make political decisions without any expertise, and the verb ἐπιψηφίζειν precisely describes the democratic process of putting a question to the vote in a democratic assembly (see e.g. *Thucydides II 24*). But, if we go back to the text, we can observe that in *Gorgias*

those instances of the lemma ἐπιψηφίζω do not intend to give criticism of the democratic process; rather, they are a reference to the fact that Socrates does know how to put a question to a vote because “there is also one whose vote I know how to take, whilst to the multitude I have not a word to say” (474a). Hence, the choice of ἐπιψηφίζειν in *Gorgias* seems coherent with the project to show the opposition between the rhetorical discourse addressed to the crowd in order to obtain the vote in the Assembly and the philosophical dialog with a sole person in order to decide if the thing in question is true or false. Thus, even if *Gorgias* is full of the insinuation that democracy ruined Athens, those mentions entail, not a real critique of democracy, but an opposition between philosophical and democratic discourse.

Another word form occurring only in *Gorgias* is ὀψοποιική (scil. τέχνη, the art of cookery), which is not a common word in ancient Greek. If we consider all the forms of this word,³⁴ Plato uses it once in the *Symposium* (187e), when the physicist Eryximachus pronounces that his craft “set high importance on a right use of the appetite for the dainties of the table, that we may cull the pleasure without the disease.” But in *Gorgias* the word occurs 7 times, mainly in the well-known passages where he makes a comparison between rhetoric and the art of cookery, precisely to show that neither are arts, but kinds of flattery (Figure 21).

³⁴ Query: ὀψοποιικ.*

tlg0059_tlg023_perseus-grc2 - 20

<465>

εἶναι τὸ τοιοῦτον, ὡ Πῶλε — τοῦτο γὰρ πρὸς σέ λέγω — ὅτι τοῦ ἠδέος στοχάζεται ἄνευ τοῦ βελτίστου · τέχνην δὲ αὐτὴν οὐ φημι εἶναι ἀλλ' ἐμπειρίαν, ὅτι οὐκ ἔχει λόγον οὐδένα ὃ προσφέρει ἢ προσφέρει ὅποι' ἄττα τὴν φύσιν ἐστίν, ὥστε τὴν αἰτίαν ἐκάστου μὴ ἔχειν εἰπεῖν. ἐγὼ δὲ τέχνην οὐ καλῶ ὃ ἂν ἡ ἄλογον πράγμα · τούτων δὲ πέρι εἰ ἀμφισβητεῖς, ἐθέλω ὑποσχεῖν λόγον. τῆ μὲν οὖν ἱατρικῆ, ὡσπερ λέγω, ἡ ὀψοποικὴ κολακεία ὑπόκειται · τῆ δὲ γυμναστικῆ κατὰ τὸν αὐτὸν τρόπον τοῦτον ἡ κομμωτικὴ, κακοῦργός τε καὶ ἀπατηλὴ καὶ ἀγεννῆς καὶ ἀνελεύθερος, σήμασιν καὶ χρώμασιν καὶ λειότητι καὶ ἐσθῆσιν ἀπατώσα, ὥστε ποιεῖν ἀλλότριον κάλλος ἐφελομένους τοῦ οἰκείου τοῦ διὰ τῆς γυμναστικῆς ἀμελεῖν. ἴν' οὖν μὴ μακρολογῶ, ἐθέλω σοι εἰπεῖν ὡσπερ οἱ γεωμέτραι — ἤδη γὰρ ἂν ἴσως ἀκολουθήσαις — ὅτι ὁ κομμωτικὴ πρὸς γυμναστικὴν, τοῦτο σοφιστικὴ πρὸς νομοθετικὴν, καὶ ὅτι ὁ ὀψοποικὴ πρὸς ἱατρικὴν, τοῦτο ρητορικὴ πρὸς δικαιοσύνην. ὅπερ μέντοι λέγω, διέστηκε μὲν οὕτω φύσει, ἅτε δ' ἐγγυς ὄντων φύρονται ἐν τῷ αὐτῷ καὶ περὶ ταῦτα σοφιστὰι καὶ ῥήτορες, καὶ οὐκ ἔχουσιν ὅτι χρῆσονται οὔτε αὐτοὶ ἑαυτοῖς οὔτε οἱ ἄλλοι ἄνθρωποι τούτοις, καὶ γὰρ ἂν, εἰ μὴ ἡ ψυχὴ τῶ σώματι ἐπεστάται, ἀλλ' αὐτὸ αὐτῷ, καὶ μὴ ὑπὸ ταύτης κατεθεωρεῖτο καὶ διεκρίνετο ἡ τε ὀψοποικὴ καὶ ἡ ἱατρικὴ, ἀλλ' αὐτὸ τὸ σῶμα ἔκρινε σταθμώμενον ταῖς χάρισι ταῖς πρὸς αὐτό, τὸ τοῦ Ἀναξαγόρου ἂν πολὺ ἦν, ὡ φίλε Πῶλε — σὺ γὰρ τούτων ἐμπειρός — ὅμοι ἂν πάντα χρήματα ἐφύρετο ἐν τῷ αὐτῷ, ἀκρίτων ὄντων τῶν τε ἱατρικῶν καὶ ὑγιεινῶν καὶ ὀψοποικῶν. ὁ μὲν οὖν ἐγὼ φημι τὴν ρητορικὴν εἶναι, ἀκήκοας · ἀντίστροφον ὀψοποιίας ἐν ψυχῇ, ὡς ἐκείνο ἐν σώματι. ἴσως μὲν οὖν ἀποπον πεποίηκα, ὅτι σε οὐκ ἔων μακροῦς λόγους λέγειν αὐτὸς συχνὸν λόγον ἀποτέτακα. ὄξιον μὲν οὖν ἐμοὶ συγγνῶμην ἔχειν ἐστίν · λέγοντος γὰρ μου βραχέα οὐκ ἐμάνθανες, οὐδὲ χρῆσθαι τῆ ἀποκρίσει ἦν σοι ἀπεκρινάμην οὐδὲν οἷός τ' ἦσθα, ἀλλ' ἔδου διηγῆσέω.

ΣΩ.

εἰ μὲν οὖν καὶ

default

PLATO170720C:"ὀψοποικ.*"

Query: Keyword: word Edit Search

sort keys: #1 None #2 None #3 None #4 None Sort

1 - 8 / 8

ref	Left context	Keyword	Right context
Symposium, s. 0187e	τῆ ἡμετέρα τέχνη μέγα ἔργον ταῖς περὶ τὴν	ὀψοποικὴν	τέχνην ἐπιθυμίαις καλῶς χρῆσθαι, ὥστ' ἄνευ νόσου τὴν ἡδονὴν καρπώσασθα.
Gorgias, s. 0463b	ἄλλα μόρια εἶναι, ἐν δὲ καὶ ἡ	ὀψοποικὴ	· ὃ δοκεῖ μὲν εἶναι τέχνη, ὡς δὲ ὁ ἐμὸς λόγος
Gorgias, s. 0464d	εἶναι. ὑπὸ μὲν οὖν τὴν ἱατρικὴν ἡ	ὀψοποικὴ	ὑποδεδυκε, καὶ προσποιεῖται τὰ βέλτιστα σιτία τῷ σώματι εἰδέναί, ὥστ'
Gorgias, s. 0465b	μὲν οὖν ἱατρικῆ, ὡσπερ λέγω, ἡ	ὀψοποικὴ	κολακεία ὑπόκειται · τῆ δὲ γυμναστικῆ κατὰ τὸν αὐτὸν τρόπον τούτων ἡ
Gorgias, s. 0465c	τοῦτο σοφιστικὴ πρὸς νομοθετικὴν, καὶ ὅτι ὁ	ὀψοποικὴ	πρὸς ἱατρικὴν, τοῦτο ρητορικὴ πρὸς δικαιοσύνην. ὅπερ μέντοι λέγω,
Gorgias, s. 0465d	μὴ ὑπὸ ταύτης κατεθεωρεῖτο καὶ διεκρίνετο ἡ τε	ὀψοποικὴ	καὶ ἡ ἱατρικὴ, ἀλλ' αὐτὸ τὸ σῶμα ἔκρινε σταθμώμενον ταῖς χάρισι
Gorgias, s. 0465d	ἀκρίτων ὄντων τῶν τε ἱατρικῶν καὶ ὑγιεινῶν καὶ	ὀψοποικῶν	· ὁ μὲν οὖν ἐγὼ φημι τὴν ρητορικὴν εἶναι, ἀκήκοας ·
Gorgias, s. 0500e	λέγειν. ἔλεγον δὲ που ὅτι ἡ μὲν	ὀψοποικὴ	οὐ μοι δοκεῖ τέχνη εἶναι ἀλλ' ἐμπειρία, ἡ δ' ἱατρικὴ,

Figure 21. Concordance of ὀψοποικὴ in the Plato corpus.

This specific use of ὀψοποικὴ has to be compared with the more common ὀψοποιία and ὀψοποιός, which appear in *Euthydemus* (290b), *Republic* (373c), and *Theaetetus* (178d), but mainly, once again, in *Gorgias* (6 instances, with 3 forms of ὀψοποιός). The frequency of ὀψοποικὴ in *Gorgias* has to be linked with the strategy already mentioned for the word “rhetoric” and more generally for words ending with -ικὴ (section 3.4). The TLG gives Plato (*Gorgias* and *Symposium*) and Xenophon’s *Oeconomicus* as the first instances of the term in Greek. According to the scholarship, the date of composition for Plato’s *Gorgias* is 385–380 (Marchand and Ponchon 2016:19), and for the *Oeconomicus* 362; there is also relative agreement that the *Symposium* is later than *Gorgias*. These claims are obviously mere hypotheses because of the lack of textual and external information about the date of composition for those texts. However, it seems that the specificity in *Gorgias* of ὀψοποικὴ, which appears when Socrates is giving a definition of rhetoric, could be an argument in favor of Schiappa’s hypothesis on the platonic origin of rhetoric. It would give also a humorous tone to the definition, if it is true that Socrates compares the emphatic new word to describe the “art” of Gorgias (so-called “rhetoric”) with an emphatic new word to describe the “art” of a cook (463b)!

3.6 Overall contrastive analysis of a corpus: Identifying the main dimensions structuring a corpus

Last but not least, for an overall view of the corpus, correspondence analysis is a multidimensional statistical tool, which computes the main dimensions of contrast structuring a corpus (see technical appendix 5.2).

We have generated such an analysis on the Plato corpus. We focused on the 200 most frequent words, which are mainly grammatical words. Each text is represented by the frequencies of its use of these words, which makes a kind of grammatical or stylistic profile (few lexical words). As explained in section 5.2, the x-axis and y-axis represent complex quantitative “mixtures” of words, and they are used to select the best 2D-perspective in the geometric representation of the data. The map (Figure 22) illustrates the relative positions of texts among one another, the main structural associations and oppositions they draw.

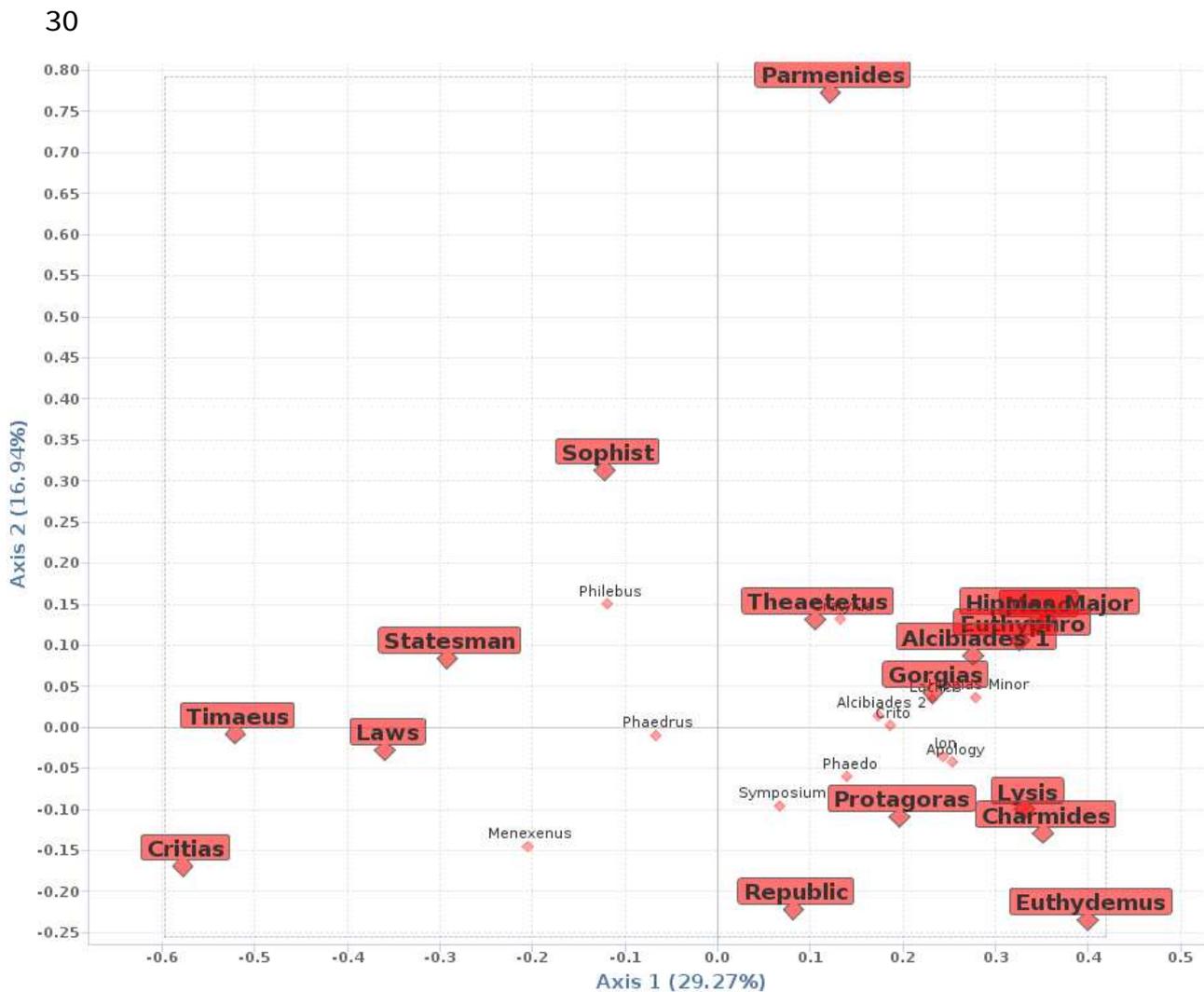


Figure 22. Correspondence analysis map of the 29 texts of Plato characterized by the 200 most frequent words of the corpus.³⁵

The main opposition we observe in the corpus is a contrast between, on the right-hand side of the map, dialog and interaction markers (ἔφη, ὦ, ἐγώ, Σώκρατες, ὅτι, μοι, σοι, εἰ, ἀλλά, ἔγωγε, οὐ, ὥσπερ, γάρ, πάνυ) and, on the left-hand side, grammatical words more associated with nouns, and then descriptive, written style (δὲ, τῶν, κατὰ, δὴ, τῆς, τὴν, εἰς, τὰς, μὲν, ὅσα, τοῖς). This written style is typical of *Laws*, *Timaeus*, *Critias*, and *Statesman*. Texts are more heterogeneous as concerns dialog and interaction; the main texts concerned are *Euthydemus*, *Gorgias*, *Meno*, *Hippias Major*, *Charmides*, *Lysis*, *Alcibiades 1*, *Protagoras*, and *Euthyphro*. Other texts are less concerned with this opposition; they may mix markers or use few of them.

Other correspondence analyses could produce different insights into the corpus: for instance, we could focus on a lexical characterization instead of a grammatical one, and try to

³⁵ Red labels emphasize texts that are best represented on this map (their quality indicator, computed as the sum of squared cosine for axes 1 and 2, is greater than 0.25).

produce a map more related to semantic and main topics. If we had morphosyntactic information, we could choose to rely on adjectives, or on verbs, as new prisms through which the corpus can be split. We could also study our corpus from the point of view of a given lexical field and its implementation throughout the texts.

4. Conclusion

4.1 Textometry's relevance for research in the Humanities

Textometry (Lebart, Salem, and Berry 1998) can be seen as a relevant approach for Humanities research, as analysis is directly grounded in texts, not necessarily mediated by external linguistic or semantic resources, which may not be fully appropriate for the corpus and thus could bias what can be seen and found. Information is, above all, taken directly from the corpus, texts, and contexts, and does not depend on available resources, which might not be precisely adjusted to every kind of data. Lexicographical resources can be used anyway, and can add valuable information, but they do not play a critical role; one does not rely on them to access and analyze the texts. The corpus can be seen, not only through an external lexicon, but also through the words it actually contains and the way they are used. This is the first important reason to be interested in textometry: texts and corpus come first, and they are a core foundation during the whole process. The corpus is the main source of information about the texts it includes, and the way those texts use words. Moreover, textometry takes the corpus as it is: one does not have to remove unknown words; one is not compelled to normalize word forms.

A second feature of textometry is that it is not an automatic approach. In a textometric analysis, the computer does not replace the researcher, it does only what it is best at: storing and processing—that's all. Of course, processing can be complex, and it often goes much further than simple counting: as shown above, textometric functions implement algorithms, score tests, and perform complex statistical functions and elaborate visualization processes.

Actually, computers are dedicated to such intensive calculations, the very kinds of tasks for which textometry uses them. Even if a computer does a great job due to its memory and processing ability, even if it allows statistical tests or analytic procedures that would otherwise be impossible, the researcher controls all aspects of the investigations, and bears ultimate responsibility for the relevance and interpretation of computed results.

While the corpus determines the context of observations and acts as a reference concerning word frequencies, its composition significantly proceeds from a human hand and a scientific

decision, not from an automatic production. This sheds light on the *quality* of the corpus (the corpus elicits interest, the researcher is familiar with it, it is carefully composed) instead of on its *quantity* (of course, a small corpus may not need any computational tool, but the hugest return of harvested data does not give any guarantee of meaningful results either).

Once the corpus has been defined, the researcher keeps on driving the analysis, finding entrance points and giving sense to the results produced by queries and statistics. Textometry is based on a strong methodology; it offers a selective choice of tools that are most relevant for textual data, but there is no unique or predefined path to get results, no unique reading of a corpus. One cannot say, “let's just try what textometry tells us about this corpus,” because there is not a given output for a given input, and even the input is a considered matter.

As shown in our examples, textometry combines quantitative and qualitative processing. Statistics or basic heuristics are defined that make sense on textual data (see section 5, technical appendix). Even qualitative KWIC concordance view is specifically designed to emphasize linguistic properties of the data. KWIC view is not just an exhaustive compilation of extracted contexts containing a given keyword. A textometric concordance combines several kinds of contextual displays: a tabulated view with context sorting, so as to reveal close contextual patterns, associated with references indicating useful upper-level metadata; and a page view in a full edition, providing all the graphical and typographical clues for reading and interpretation.

In such an analytic framework, textual editions can still meet requirements for textual studies in the Humanities, according to the researcher's experience that reading, analyzing, and editing, are not separate activities. TXM software, for example, can work on TEI encoded corpora, display fine HTML editions, and manage several aligned versions of the same text, including multimedia capture of the original source (for instance pictures of manuscripts).

4.2 Experiencing textometry on new corpora

The purpose of this paper was not primarily to establish new scientific results about Plato's *Gorgias*, but it aims to introduce the classicist scientific community to the textometric approach to digital studies. Textometry could be characterized as a both qualitative and quantitative analysis methodology. Behind a detailed typology of queries that could be searched for in the corpus (queries about word frequency, syntagmatic or paradigmatic uses of word meaning in the corpus, word evolution, texts' contrastive characterization and corpus internal structures),

the key idea is that it is possible in a digital context to elaborate tools combining advanced computational analysis and philologic attention and sensibility towards text integrity and richness.

We hope that this approach can become a personal experience for many colleagues, since several applications implement the textometric ideas promoted here. Our example has been realized with TXM open-source software, which is available for multiple operating systems (including a web portal version), most languages (including Latin and ancient Greek), and many corpus encoding states (from raw text corpora to TEI encoded ones). It also gives all the resources to put into TXM any corpus built from the open Perseus Digital Library.

Acknowledgments

We are grateful to our two reviewers for their accurate comments and stimulating suggestions.

5. Technical Appendix: More Information about Three Textometric Processes

5.1 Specificities: A statistical computation to find keywords

The specificity score (Lebart, Salem, and Berry 1998) is an application of Fisher's exact test to textual data (Gries 2012). Specificity analysis allows one to identify which words are specifically overused and underused in a part of the corpus, compared to its use in the whole corpus. It is like a keyword analysis in some textual analysis programs, but instead of using a log-likelihood, t-score, z-score, or tf.idf measure, it implements Fisher's exact test, which is known as the most accurate measure for words' frequency distributions (Gries 2014; McEnery and Hardie 2012).

Based on a hypergeometric probability model, this analysis provides either a positive or negative specificity score. A *positive* score indicates the order of magnitude of the probability of a word w appearing f times *or more* than f times in a part containing n words, given that w appears F times in whole corpus of N words. A *negative* score is given when the frequency is lower than expected on a random basis; the measured probability is then the one for a word appearing f times *or less* than f times, given its total frequency F , and part and corpus sizes (n and N). Specificity scores equal to or higher than 3 (1 chance in 1,000 to obtain the frequency f or more if the words were randomly distributed among parts) or lower than -3 (1 chance in 1,000 to obtain the frequency f or less randomly) are considered significant.

Figures 6, 7, 9, and 20 show examples of specificities outputs, which point out lexical choices of some texts within Plato's work.

5.2 Correspondence analysis: A geometrical optimization to get a word-based map of the corpus

Map visualization of the corpus is based on correspondence analysis, a multivariate statistical tool. The idea behind correspondence analysis may be explained in a few words. Every text is represented according to the words it uses: geometrically speaking, each text is a point and gets a unique location in a multidimensional space in which each word is an axis. The coordinate of a text point on the w word axis is related to the frequency of the w word in the text. If two texts have many words in common and use them with similar frequencies, then their points get close locations in space. Mathematical transformations are then used to build a new set of axes for this space, with new worthwhile properties. The new axes combine the previous word-axes (each axis is a linear combination of word-axes) so that: (i) all the information is kept (points representing texts keep their distances to one another, the shape of the "text-cloud" is the same); (ii) redundancy is eliminated, which reduces the number of axes; (iii) the new axes are ordered so that the first one shows the "largest" and "heaviest" variations (i.e. along this axis frequencies are very different and concern many words and texts), and each following axis brings the "largest" and "heaviest" remaining variations, so that with axes from 1 to any number N , one gets the best N - dimensional view of the lexical variations and text oppositions inside the corpus. Given the value 2 for N , the corpus can then be visualized as a map using the Cartesian coordinate system, knowing that this map is mathematically proven as being the best 2D-visualization in order to focus on the greatest differences inside the corpus. What is to be read on this map is a quantitative summary of the main correspondences and oppositions among texts. The axes are complex weighted combinations of words whose composition might be studied, but the main use of correspondence analysis for text analysis is to look at the relative positions of texts, and axes are first used to provide the best angle to get the best view.

Figure 22 applies correspondence analysis to Plato's work, so that one gets a view about text similarities or oppositions for high frequency word usage.

Correspondence analysis is conceptually similar to principal component analysis, but correspondence analysis is preferred to principal component analysis here because it better fits the data. In the field of textometry, textual data are represented in frequency tables where rows are words (or other kinds of linguistic units), and columns are texts (or groups or parts of texts

dividing the corpus). In this kind of table, rows and columns are both categorical variables (a set of words and a set of texts) and play symmetric roles, so structurally this is a two-way table (also called a cross tabulation, or a contingency table). Correspondence analysis is specifically designed and relevant for the analysis of such contingency tables, for which it proves getting better results than principal component analysis (Lebart, Salem, and Berry 1998:63-69).

5.3 “Back-to-text”: Concordances and word-in-context functions as main features

What is called the “back-to-text” functionality in the textometric field is actually the core of textometric processing: every result (lexical list, table, or graph) must be interpreted with an eye toward corresponding words in context.

This could be a distinctive feature of a textometric approach, an efficient way to distinguish it among many current text mining proposals. With a text mining approach, one digs into texts in order to extract pieces of information, so that all the analysis focuses on these few well identified pieces instead of an unstructured full text; and visualization is often the goal, the end of the analysis, with a synthetic and suggestive view that replaces and summarizes the corpus. On the contrary, the textometric approach chooses and implements the hermeneutic circle: any statistical summary or view has to be interpreted looking back to the text, any distant reading view invites renewed close reading, so analysis cannot get rid of textual richness and complexity—it is grounded on this textual source material.

Bibliography

- Benzécri, Jean-Paul, *et al.* 1973. *L'analyse des données*. Vol. 1, *La taxinomie*. Vol. 2, *L'Analyse des Correspondances*. Paris.
- . 1981. *Pratique de l'analyse des données*. Vol. 3, *Linguistique & lexicologie*. Paris.
- Brunet, Étienne. 2009. *Écrits choisis*. Vol. 1, *Comptes d'auteurs: Études statistiques de Rabelais à Gracq*. Ed. Damon Mayaffre. Paris.
- . 2011. *Écrits choisis*. Vol. 2, *Ce qui compte: Méthodes statistiques*. Ed. Céline Poudat. Paris.
- . 2016. *Écrits choisis*. Vol. 3, *Tous comptes faits: Questions linguistiques*. Ed. Bénédicte Pincemin. Paris.
- Christ, Oliver. 1994. “A Modular and Flexible Architecture for an Integrated Corpus Query System.” In *Papers in Computational Lexicography (Complex '94)*, 22–32. Budapest.

- Evert, Stefan, and Andrew Hardie. 2011. "Twenty-First Century Corpus Workbench: Updating a Query Architecture for the New Millennium." In *Proceedings of the Corpus Linguistics 2011 Conference*. Birmingham.
- Gries, Stefan Th. 2012. "Corpus Linguistics: Quantitative Methods." In *The Encyclopedia of Applied Linguistics*, ed. C. A. Chapelle, 1380–1385. Oxford.
- . 2014. "Quantitative Corpus Approaches to Linguistic Analysis: Seven or Eight Levels of Resolution and the Lessons They Teach Us." In *Developments in English: Expanding Electronic Evidence*, ed. M. Kytö Taavitsainen, Cl. Claridge, and J. Smith, 29–47. Cambridge.
- Heiden, Serge. 2010. "The TXM Platform: Building Open-Source Textual Analysis Software Compatible with the TEI Encoding Scheme." In *24th Pacific Asia Conference on Language, Information and Computation - PACLIC24*, ed. R. Otoguro, K. Ishikawa, H. Umemoto, K. Yoshimoto, and Y. Harada, 389–398. Tokyo.
- Heiden, Serge, Matthieu Decorde, and Sébastien Jacquot. 2018. *TXM User Manual: Version 0.7 Alpha*. Trans. Sara Pullin. Lyon. <http://textometrie.ens-lyon.fr/files/documentation/TXM%20Manual%200.7.pdf>.
- Lafon, Pierre. 1984. *Dépouillements et statistiques en lexicométrie*. Paris.
- Lebart, Ludovic, André Salem, and Lisette Berry. 1998. *Exploring Textual Data*. Boston.
- Léon, Jacqueline, and Sylvain Loiseau, eds. 2016. *Quantitative Linguistics in France*. Lüdenscheid.
- Marchand, Stéphane, and Pierre Ponchon. 2016. *Gorgias de Platon suivi de Éloge d'Hélène de Gorgias*. Paris.
- Marchello-Nizia, Christiane, and Alexei Lavrentiev. 2013. *Queste del saint Graal: Édition numérique interactive du manuscrit de Lyon (Bibliothèque municipale, P.A. 77)*. Lyon. http://catalog.bfm-corpus.org/qgraal_cm.
- McEnery, Tony, and Andrew Hardie. 2012. *Corpus Linguistics: Method, Theory and Practice*. Cambridge, MA.
- O'Sullivan, Neil. 1993. "Plato and ἡ Καλουμένη Ῥητορικὴ." *Mnemosyne* 46:87–89.
- Salem, André. 1987. *Pratique des segments répétés: essai de statistique textuelle*. Paris.
- Schiappa, Edward. 1990. "Did Plato Coin *Rhètorikè*?" *The American Journal of Philology* 111:457–470. <https://doi.org/10.2307/295241>.
- . 1999. *The Beginnings of Rhetorical Theory in Classical Greece*. New Haven and London.