

# Distant Philosophizing

## Applying computational literary studies to philosophical issues

Alessandro Rosa

June 2021

## 1 Introduction

For centuries philosophers have debated theoretical concepts related to the true essence of reality. This includes discussing the very meaning of things such as Truth, Beauty, Justice or to inquire into the elementary features of Nature itself. Many of these issues were raised at the very beginning of philosophy as an intellectual practice, namely with Greek thinkers such as Thales of Miletus, Socrates and Plato. They posed many fundamental questions and provided some answers, many of them still debated at the present time.

Most intellectuals across the cultural development of Western thought tried to solve those issues, leveraging on the cultural context in which they lived. Whoever has worked on a philosophical essay knows that it is almost impossible to avoid referring to past philosophers that have, in some way or another, tried to tackle the very same questions. But historians of philosophy have extensively discussed whether there is a fundamental difference in how those issues were raised (and solved) in the early years of philosophy and how they are currently debated. The question can be posed as follows: are we discussing the same things that Plato discussed during his time? The question, of course, is related with another fundamental question: can philosophy make any progress, in the same sense in which physics or biology do?

The aim of this work is to exploit computational tools and quantitative methods to verify whether such research hypothesis could be validated or not with empirical evidence. It is an attempt to apply the Distant Reading methodology that is currently employed in critical literary studies to philosophical problems. It is not an attempt to answer in depth those philosophical questions (although it might provide some help for those who want to do so). In order to do so, a little experiment involving a specific philosophical issue, namely the problem of the nature of knowledge, has been carried out using a custom implementation of Python packages. The results and the whole code can be found at the [GitHub repository](#) of the project.

## 2 Distant Reading: a methodology

The notion of Distant Reading was firstly introduced by Franco Moretti at the beginning of the 2000s, and rapidly became popular in the scholarly debate. The basic idea behind Distant Reading is that we might want to distance ourselves from the literary works under examination in order to have a better understanding of the macro phenomena that are in play. By exploiting computational tools we are able to perform a different kind of analysis of literary texts, one than involves probabilistic inferences and statistical reasoning. For scholars involved in the literary field this is a very distant activity from the so-called Close Reading, i.e. the traditional approach that consists in reading in depth a piece of literature to extract information and provide conjectures about cultural phenomena. This is particularly true for problems that spread across time that are hardly detectable using only a closer look, and for which a wider perspective is fundamental.

Let me consider a paradigmatic example provided by Ted Underwood [1]. Both scholars and every-day readers are familiar with the concept of literary genre. We are usually able to distinguish between different kind of fictional novels (e.g. detective stories, science fiction, horror, romantic novels and so on) because we can detect some peculiar features of a story (for instance, the book might open with a crime scene) and make inferences regarding its genre. However, the boundaries between genres are not always that sharp, and scholars might want to argue in depth what those peculiar features are and in which sense two books belong to different genres. They might want to explain why those books diverge and how a specific literary genre has evolved through the centuries. Moreover, they might want to answer those questions without relying entirely on their own personal perspectives.

As a consequence, scholars who focused on these kind of questions could find computational and quantitative methods rather useful for their scopes. Distant Reading enables scholars to look at a larger picture, taking into account more data and evidence to validate their research hypothesis. It allows to explain through empirical evidence macro phenomena related to literature that otherwise would not even be possible to see. Humanists are usually reluctant to use quantitative methods and employing statistical analysis for their researchers, but there is little doubt that this could provide more scientific support to the field. This is not to say that qualitative studies have no room within the Distant Reading approach. One important thing to be aware of, and which is often reminded by Underwood himself, is the fact that "numbers are not inherently more or less objective than words". The goal is not to completely get rid of qualitative studies and personal beliefs, but to provide tools that dispense empirical evidence and enable us to discover new interesting things about the subject. And since personal perspective is not completely eliminated by the process, it should also be reminded that quantitative studies might include biases of any sort exactly as the traditional approach does.

One important feature of Distant Reading is the fact that we can analyse thousands of works in few minutes without actually reading them. Since scholars

tend more and more to have less time to spend on their own research, there is an indubitable advantage in doing so. But despite the pragmatic benefits, being able to read a potentially infinite number of texts in short time allows us to include in our studies also less known works that were almost forgotten. There is an endless list of authors that produced many literary works, being them novels, poetry or essays, that have never been recognized as important enough to be included in literary studies and ended up in what Moretti called "The Slaughterhouse of Literature" [2]. Italian high school students are very familiar with Dante, Petrarca and Cavalcanti, but what about all the other poets of that time that were overshadowed by talent of the *Commedia*'s author? We might want to stop associating our concept of literature with names of famous and well known authors and starting considering literature as a whole, including unrecognized authors from different countries, speaking different languages and belonging to different cultures and subcultures.

### 3 Philosophizing: the case study

Philosophy is renowned to have raised conceptual problems for centuries. The philosophical tradition includes an endless list of authors and intellectuals that started to question and analyze in depth the very essence of reality since the VI Century BC. Despite being an old tradition, many (if not all) philosophical problems are still unsolved and still debated in today's academical classes. Many of these problems were faced multiple times in different periods and in different cultural contexts. The very notion of what it means for something to be true, for instance, was first questioned by Plato and by his most famous student Aristotle, then it was discussed in medieval times within the Catholic religious context, then again questioned during Humanism and it is hard to claim that nowadays there is a widely accepted view on the subject.

A common issue for scholars that are interested in the evolution of philosophical thought is to understand whether the subject has undergone a proper development throughout the centuries or not. In other words, they are interested in knowing if the philosophical problems (and their possible solutions) that were first raised by ancient philosophers are the same that we are discussing nowadays. The issue could be posed as the following: when we now inquire into the very meaning of the concept of Truth, is it the same "thing" that was initially conceptualized by Plato? For instance, we could easily imagine that for medieval authors Truth was very close, and in some important sense influenced by, the concept of God. Indeed, deities played an important role for Plato but we know that the Greek approach to religion was rather different from Christianity. However, both Plato and (let us say) Augustine debated around the same thing, i.e. what do we mean when we talk about Truth.

Despite major differences in their cultural contexts it could be argued that there is little if no difference between Plato's conceptualization of Truth and Augustine's. Alfred North Whitehead notoriously claimed that it is possible to reduce all western philosophy to the problems that were raised by Plato

i.e. that in his writing there already were all the possible ideas that we could conceptualize about things like Truth, Beauty, Nature, the human Soul and so on. The core idea behind this paper is that we might want to argue against such claims, and to do so Distant Reading could be particularly handfuf.

For the present experiment I propose to analyse a specific problem, namely the problem of what can be accounted as knowledge. The choice is due to the fact that it is a specific but wide issue that has been discussed fiercely throughout all the two thousand years of western philosophy. Moreover, it is possible to precisely identify philosophical essays that tackled the topic since the very beginning i.e. since Plato's dialogues (in fact, it is Plato himself to raise the issue, as it is commonly the case).

The problem of knowledge can be summarized as the following: how can we claim that we know a specific thing? As many philosophical problems, the answer seems pretty straightforward: we know something because we have experience of it (we see/touch/hear it). However, reducing knowledge to experience seems very limiting because we feel that there is a fundamental difference between the two. Thus, philosophers want to inquire into this apparent difference. Plato's solution involves a three-part requirement: we have knowledge if we have a justified opinion that turns out to be true about something.

This is not the place to discuss such issue in depth, so I will just notice that all of the three requirements are potentially problematic per se (who decided what is "justified"? What is a belief? How can we know that something is true? And how is truth related to knowledge?). The interesting thing for us is to understand whether the problem of knowledge was tackled in different ways by different authors without actually reading all the philosophical essays that were written on the subject. If, as I am going to argue, we could do that, it would become possible to gain interesting hints about how philosophers deal with philosophical issues, organizing different approaches in different clusters and maybe proposing different ways of grouping cultural movements.

## 4 Development: quantitative tools for philosophers

For this project I decided to use a limited number of philosophical works ranging from ancient Greek philosophy to contemporary analytical philosophy. I selected 16 texts, more specifically: 3 Plato's dialogues, 1 large essay by Aristotle, 2 works by Descartes, 1 by Hume, 2 by Kant, 4 by Ernest Sosa, 2 by Timothy Williamson and 1 by Richard Feldman. The works were selected for their relevance to the philosophical problem we are dealing with, i.e. with the problem of knowledge. Of course, the selection might be argued against for its philological accurateness: The Critique of Pure Reason deals with a lot of issues, not merely the problem of knowledge, and the same can be said for most of philosophical essays. Moreover, it is clearly a small data set to work with and we should not expect statistically relevant results, which can be drawn only by

using a large data set. However, the goal here was to lay a general framework that could be expanded later in the future with additional textual files.

All the texts were extracted in a pure textual format (.txt) and, in some cases, converted from .pdf format using [pdfminer.six](#). Data was firstly analysed with the [Python nltk library](#) in order to extract the most relevant textual features, pre-processing the files in order to tokenize the data set and removing stop-words. Some additional cleaning was required in order to get rid of undesired results due to the peculiar features of some texts. For instance, all Plato's dialogues include the names of the speaker at the beginning of a sentences; thus, those names recurred very often in the text and eventually ended up being the most frequent words for Plato's works. Although the raw frequency of the number of lines said by a character may be of some relevance for other implementations (for instance if we wanted to create a network of the dialogue's protagonists), it did not make much sense to include it for the present work.

The data set was then split into three clusters, namely the ancient, modern and contemporary ones. This allowed to have a general overview of the differences between three philosophical macro periods. Both the raw word occurrences and the TF-IDF[3] scores illustrate that there is a clear semantic difference in how the same problem was tackled by those philosophers. For instance, in ancient philosophy the problem of knowledge seems to be strictly related with the concepts of "soul", "mind" and "virtue". The same cannot be said neither for modern philosophy, which shows much more interest in the concepts of "reason", "idea" and "object", nor for contemporary academics, that are now talking about "justification", "evidence" and "belief".

For a more accurate and interesting result a small topic modelling analysis was implemented. Topic modelling[4] is an algorithmic technique that tries to subdivide words inside a text (or a corpus) into clusters that are somehow semantically meaningful. The underlying assumption is that words that usually come together are associated with a specific topic. Thus, a topic is nothing but a list of words that capture a specific semantic meaning (presumably) wanted by the author. Usually topic modelling is implemented to have a general understanding of a large corpus of texts, for instance to individuate macro themes and more generally what the corpus is about.

However, here topic modelling was applied with a different scope. We already know the general topic that we are dealing with (the problem of knowledge) and we want to see whether we can recognize a variation in the semantic clusters when we divide the whole corpus into philosophical subsections. In order to do so, [Gensim](#), a Python library for topic modelling, document indexing and similarity retrieval with large corpora, was implemented. Gensim basically transforms the documents in the corpus into vectors and exploits various algorithms to perform topic modelling. For sake of simplicity, in the current project LDA (Latent Dirichlet Analysis) was implemented without any changes to the basic settings. LDA is a popular method for topic modelling which considers each document as a collection of topics in a certain proportion.

Since analyzing topic modelling results can be a very complex task, a proper visualization tool was implemented to easily draw some conclusions. [LDAvis](#)

allows to easily interpret Genism output by providing a visualization tool [5] with which users can immediately see the semantic meaning of each topic. Selecting a topic on the left allows to visualize the most useful terms for that topic, along with a global rank of the most relevant words; selecting a word reveals the conditional distribution over topics of the selected term. The area of the circles are proportional to the relative prevalence of the topics in the corpus, while the distance between the centre of each circle represents how much two topics diverge from each other. The  $\lambda$  value determines the weight given to the probability of term  $w$  under topic  $k$  relative to its lift (measuring both on the log scale). Authors suggest to keep the  $\lambda$  value at 0.6, because setting it to 1 would mean that terms would be ranked solely by their lift, which implies that the red bars would be sorted from widest (at the top) to narrowest (at the bottom). The results can be fully explored in an HTML form at the [GitHub repository](#).

## 5 Conclusion

The goal of this brief paper was to show whether it would be possible to implement computational methods in order to validate scholar hypothesis related to historiography of philosophy. The take home lesson here is that by exploiting quantitative studies it is possible to have a general look at the way in which the philosophical tradition has dealt with a problem. I decided to use a simple and general subdivision (namely a chronological one), but the idea is that by analysing a large corpus of text it might be possible to explore different way of clustering philosophical essays. More importantly, it would allow us to consider philosophers that are less well known than Plato or Kant, and whose contributions might be passed under the radar. It also makes possible to have a wider consideration of philosophy as a global activity, for instance including Asian philosophers in our studies and confront their approach to a philosophical problem with the Western one.

The results that were obtained illustrate that something like that is quite possible, but it is important to specify the main limitations that this study carries. First of all, it is clear that the number of works analysed is insufficient to obtain statistically significant results. When dealing with textual data it is fundamental to have a large corpus to explore in order to avoid random fluctuations. But there are also some important methodological issues that are specific of this case. When dealing with philosophy, it is necessary to clearly identify the boundaries of the scholar hypothesis we want to validate or confute. This means identifying a specific problem and the historical time periods that have dealt with it, the authors involved and most importantly the philosophical essays that have faced the problem. But philosophical essays have rarely tackled only one single issue, and even when they do they are usually inscribed within a conceptual framework that might be not completely devoted to the single philosophical issue we are interested in. It might be the case that we cannot make a sharp distinction between one philosophical issue and another, leading

to include "unreliable" texts in our research.

Another important issue to take into account is the problem of translation. Since most of the past philosophers have written their essays in languages different from English, the texts that were analysed here were taken from translated editions. This is particularly problematic for ancient texts, such as Greek ones, because the accuracy of translation of certain specific terms might be debatable. Scholars frequently argue for or against a specific translation, and thus it might be important to decide beforehand which critical edition to consider when feeding a text to an algorithm.

Despite these limits, I believe that such computational tools and quantitative analysis could play a fundamental role for finding new empirical evidence about research hypothesis that were traditionally only limited to the scholar's point of view on the subject. Following the leading path of literary studies, Distant Reading might prove to be a valuable resource in philosophical research, especially when dealing with wide range phenomena.

## References

- [1] Ted Underwood. *Distant Horizons: Digital Evidence and Literary Change*. University of Chicago Press, 2019.
- [2] Franco Moretti. The slaughterhouse of literature. *MLQ: Modern Language Quarterly*, 2000.
- [3] Matthew J. Lavin. Analyzing documents with tf-idf. 2019.
- [4] Shashank Kapadia. Topic modeling in python: Latent dirichlet allocation (lda): How to get started with topic modeling using lda in python. 2019.
- [5] C. Sievert and K. E. Shirley. Ldavis: A method for visualizing and interpreting topics. *Association for Computational Linguistics*, 2014.