

Exploratory Text Analytics Final Project
Analysis of Classical Greek and Roman Works

Author: Alexandra Ferentinos

DS5001

Date: 12-06-2024

Introduction:

The stories of the pantheon, the gods, heroes, monsters, and love weave literary narratives that paint the stories of the beliefs and life perspectives of the past. This report explores the corpus of Classical Greek and Roman literature, as they are often considered foundational texts that still inspire modern literature, media, and pop culture today. The corpus includes nine classic works spanning epic poetry, drama, and mythological narratives, from Homer's well-known epics to Ovid's tale of metamorphosis.

	book_id	title
doc_id		
0	31	Sophocles-31: Plays of Sophocles: Oedipus the ...
1	228	Virgil-228: The Aeneid
2	348	Hesiod-348: Hesiod, the Homeric Hymns, and Hom...
3	830	Apollonius Rhodius-830: The Argonautica
4	1727	Homer-1727: The Odyssey
5	6130	Homer-6130: The Iliad
6	7700	Aristophanes-7700: Lysistrata
7	21765	Ovid-21765: The Metamorphoses of Ovid, Books I...
8	35451	Euripides-35451: Medea of Euripides

Figure 1: DOC table displaying the Book ID and Title of the works in the corpus.

Figure 1 shows the overview of the titles of the works in the corpus. The works range from Sophocles' Oedipus trilogy that explore themes of fate vs free will; Homer's Iliad and Odyssey which describe the events of the Trojan war and Odysseus homecoming; and Ovid's Metamorphoses, which show the perspectives of the gods vs man, as well as what is morally right vs wrong. This corpus covers a robust source of universal themes that will be explored in the exploratory text analysis of these works, such as divine intervention, heroism, and free will.

Using natural language processing techniques learned over this semester mainly: topic modeling, and sentiment analysis, this report aims to analyze the patterns present in the authors approach of classical themes. The themes of divine intervention, heroism, free will, and many others.

Data Collection:

The data was collected from the Project Gutenberg website that has a collection of over 70,000 free eBooks for use and download.ⁱ The texts were downloaded as plain text UTF-8 files that contain Gutenberg metadata, as well as the classical text. Each text file was approximately over 10,000 lines long and were stored in a storage directory. The text files were then

individually analyzed by hand to determine the proper regex format required per work to discern the chapter formatting, as well as unnecessary metadata removal.

The nine texts stored in the directory were then accessed by the `acquire_epubs` function used and implemented throughout the semester. These texts were then tokenized, reduced, and a document table and a vocabulary table were created (F1-F5 Text data forms). When the Token table was made it still contained stop words; this would have added noise to the later analysis of the works. A dictionary of common stop words was made and used to filter it down to a reduced token table. This reduced table was added to the csv file read in pipeline to streamline the analysis of the works.

Analysis:

As mentioned previously, the corpus is comprised of nine classical Greco-Roman texts with the following categories and distributions of tokens: Epic Poetry category; Homer's *Iliad* and *Odyssey* (approx. 147,000 tokens), Virgil's *Aeneid* (approx. 56,000 tokens), Ovid's *Metamorphoses* (approx. 58,000 tokens), Apollonius' *Argonautica* (approx. 21,000 tokens). Drama category; Sophocles' plays (approx. 13,500 tokens), Euripides' *Medea* (approx. 7,500 tokens), Aristophanes' *Lysistrata* (approx. 6,300 tokens). Finally, the Historical category (approx. 13,000 tokens). The categories are useful for analysis and when interpreting trends in the Zipf distribution, as shown in Figure 2. When interpreting the graph, there is a concentrated density of Epic Poetry category tokens in the upper ranges of the plot, while the Drama category is densely comprised that of the mid-range of the graph.

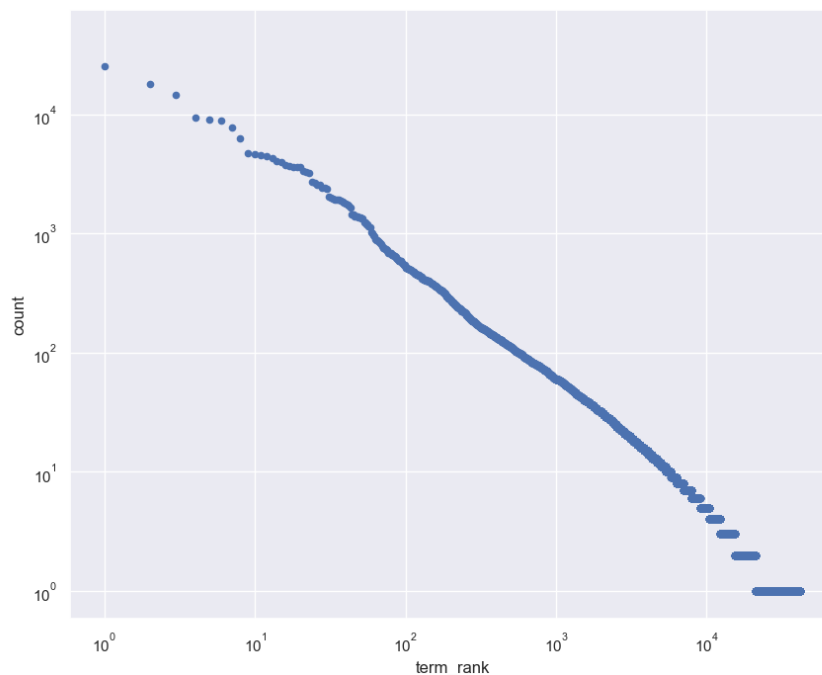


Figure 2: Work Frequency Zipf distribution of the Classical Corpus.

The TFIDF significance matrix of the classical text corpus was effective in illuminating the patterns of term usage across the various texts of the corpus. Figure 3 is a graph of the heatmap of the TFIDF weights. It incorporates term rank, polarity, and significance scores across the texts, where the darker blue color indicates a higher significance score.

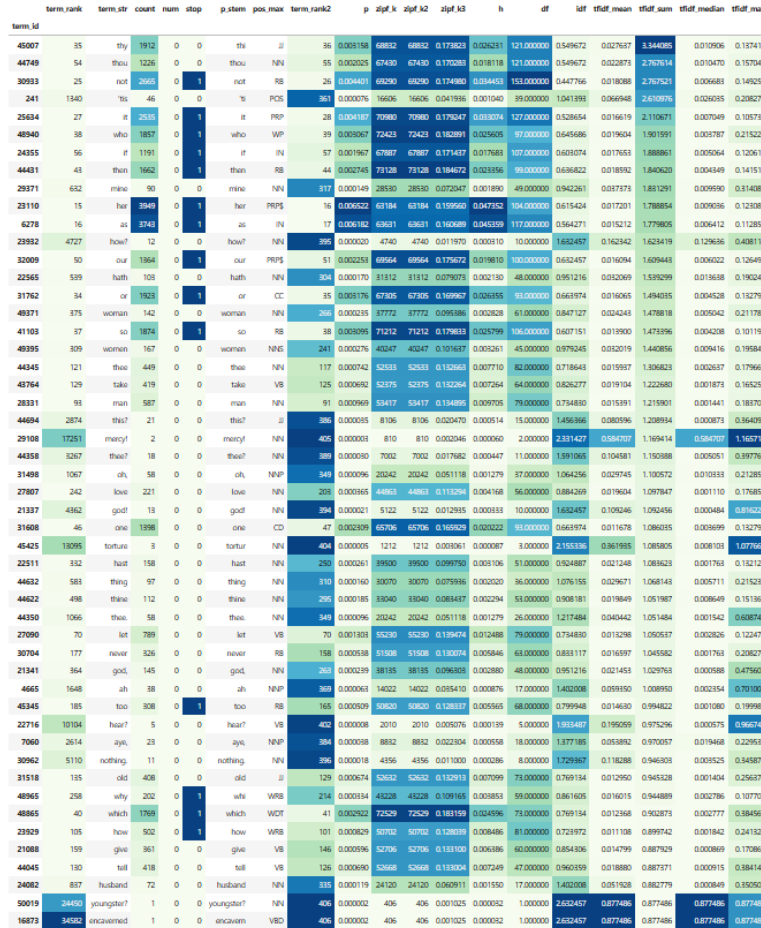


Figure 3: TFIDF Matrix Heatmap of the Classical Texts.

The main takeaways from this heatmap are that the Epic Poetry category has high frequency of terms associated with “gods”, “fate”, and “heroes.” This makes sense for this genre as it is comprised of themes such as being propositioned by the gods, and the hero’s journey. The Drama category has a high frequency of terms like: “chorus”, “revenge”, and “women.” This coincides with the classics drama genre of the chorus announcing moments of tension and women being primarily associated with acts of revenge and drama in this time. Finally, the History category has moderate frequency usage of words like “hero”, “travel”, and “journey”. While these terms can be more ambiguous, it makes sense as this is a smaller category. Overall, the heatmap demonstrates the patterns per category parallel to their category themes.

Figure 4 gives insight into the LDA analysis of the theme clusters across the Greco-Roman corpus. There appears to be three key topic clusters in terms of the content groupings. First, the brown colored cluster grouping centers around daily life, with the words being related to domesticity. However, if one knows Odysseus from the Odyssey the terms of domesticity have a much darker context in relation to his retaliation towards the suitors pursuing his wife, Penelope. Second, the green and blue areas relate to the interactions of the gods and man. Finally, the orange section relates to mentions of war and strategy.

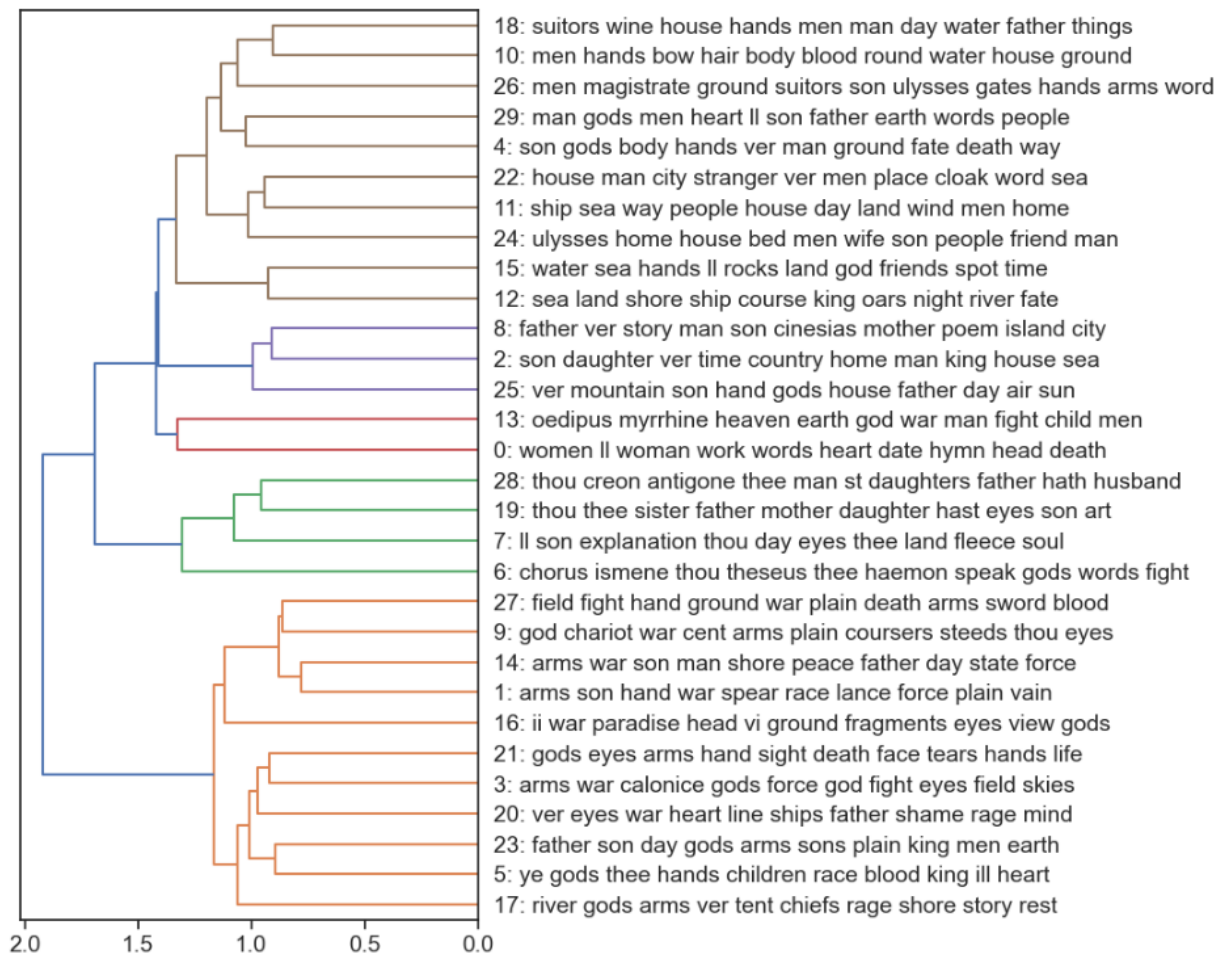


Figure 4: Hierarchical clustering dendrogram of the Classical texts.

Figure 5 shows the clustering of the principal components across the corpus. This visualization further demonstrates the thematic divide amongst the works. Texts like “The Iliad”, “The Odyssey”, and “The Aeneid”, which are all Epics, are divided from the Dramatic category works of “Medea” and “The works of Sophocles.” The others’ work has a slight overlap amongst the themes, while still preserving some of their independence. A good example of this is “The Argonautica”, as it is in between the separation of the Epic works and the Dramatic ones. This provides additional context to how there is thematic independence amongst some of the works but shared themes between others.

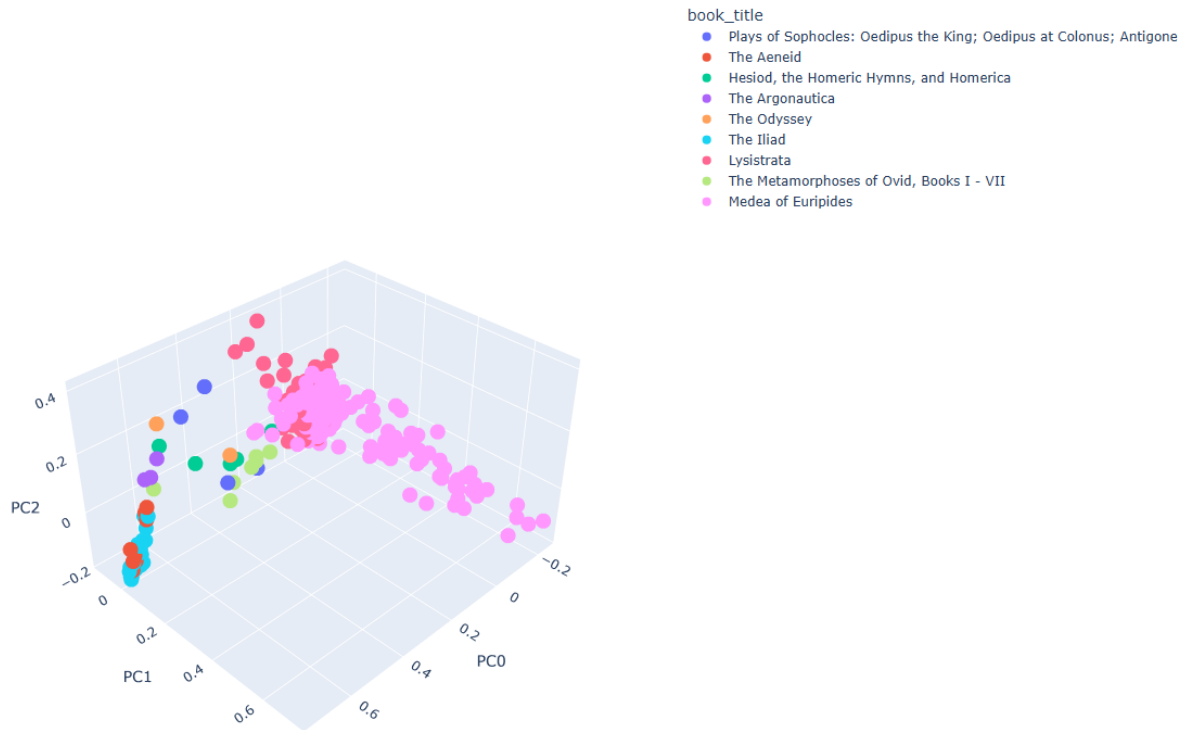


Figure 5: PCA plot of the clustering of principal components of the classical corpus, Three-dimensional format.

An overview of the word embedding semantics across the corpus text at each category level gives more insight into the patterns between the authors and category genre. Figure 6 of the Epic Poet category shows an apparent correlation between vocabulary that has to do with the divine, “gods”, and marital themes. This correlation is fitting as the works of Homer, Virgil, Ovid, and Apollonius rely heavily on narratives of the will of the gods and the effects it has on the mortals’ lives, specifically their marriages. The most notable example is Odysseus’s relationship with the will of the gods and his wife Penelope in the *Odyssey* and *Iliad* by Homer.

Epic Poets Comparison (Homer, Virgil, Ovid, Apollonius)

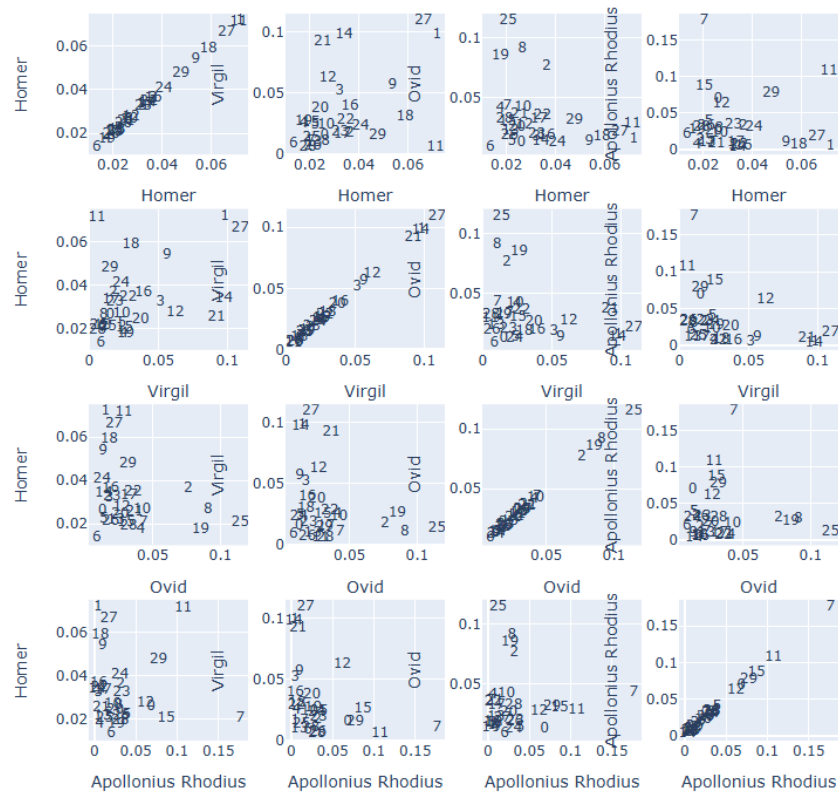
**Figure 6:** Word embedding pairwise comparison between Epic Poets.

Figure 7 shows the word embeddings comparison between the Dramatic Poets. The clustering behavior seen in the works of Euripides and Sophocles to the tragic vocabulary aligns with their works motifs of the darker tales of irony. The tale of Oedipus with its tragic ironic fate of love and lust is one example. The work of Aristophanes leans away from the tragic vocabulary into the extreme of comic humor. This behavior in the Dramatic Poets category is rational because Drama encompasses aspects of tragedy, irony, humor, and comedy. Finally, Hesiod's work moves amongst the scales of tragedy and comedy, as this text is that of a more historical perspective of the period.

Dramatic Poets and Hesiod Comparison

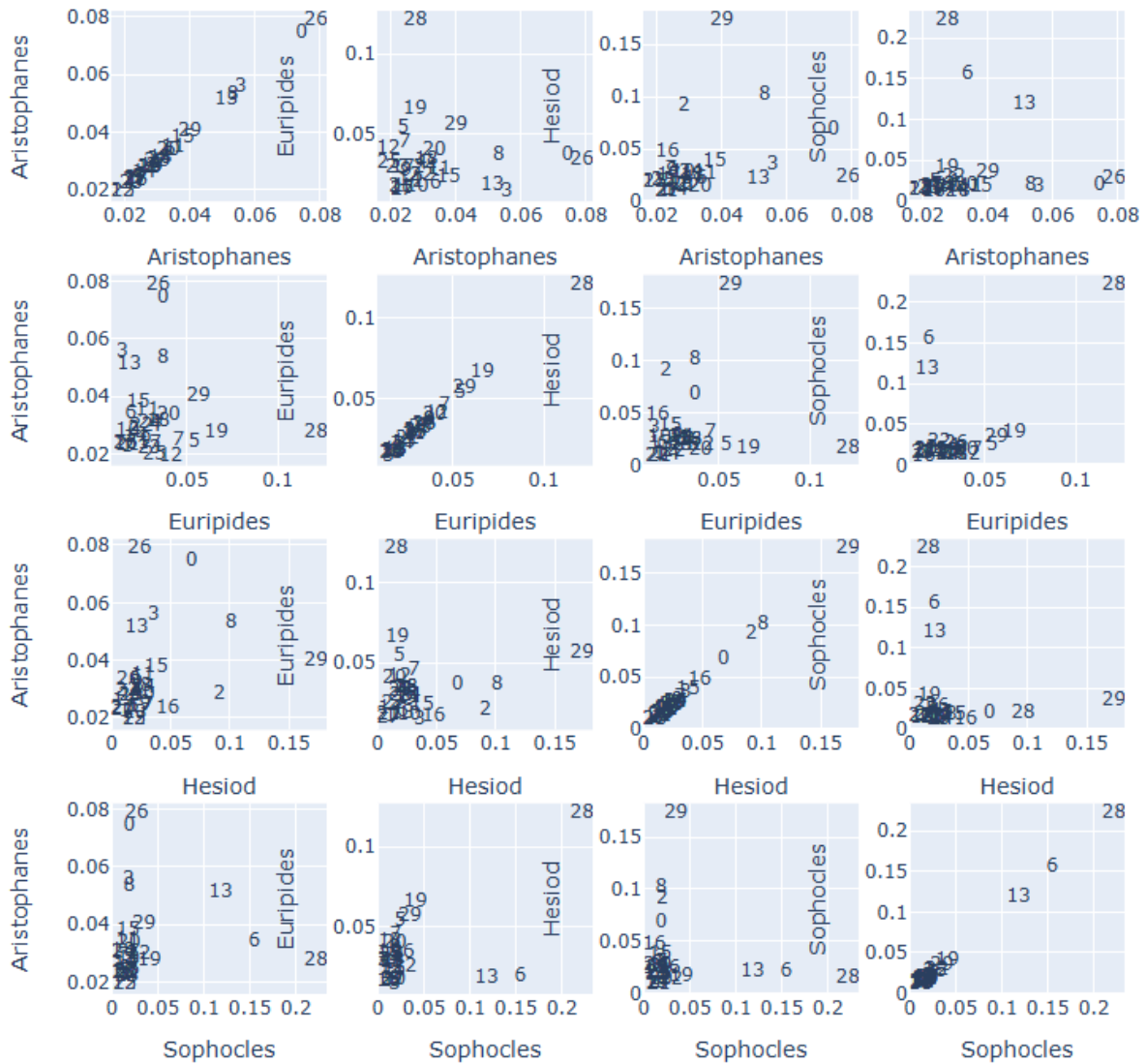


Figure 7: Word embedding pairwise comparison between Dramatic Poets.

Figure 8 is a visual representation of the word embedding pairwise comparison between Epic and Dramatic Poets. This graph shows distinct separation between the Epic and Dramatic corpus texts vocabularies. This visualization also highlights the shared vocabulary related to language of mythos, like “divine,” “nymph,” “chorus,” etc.

Epic vs Dramatic Poets Comparison

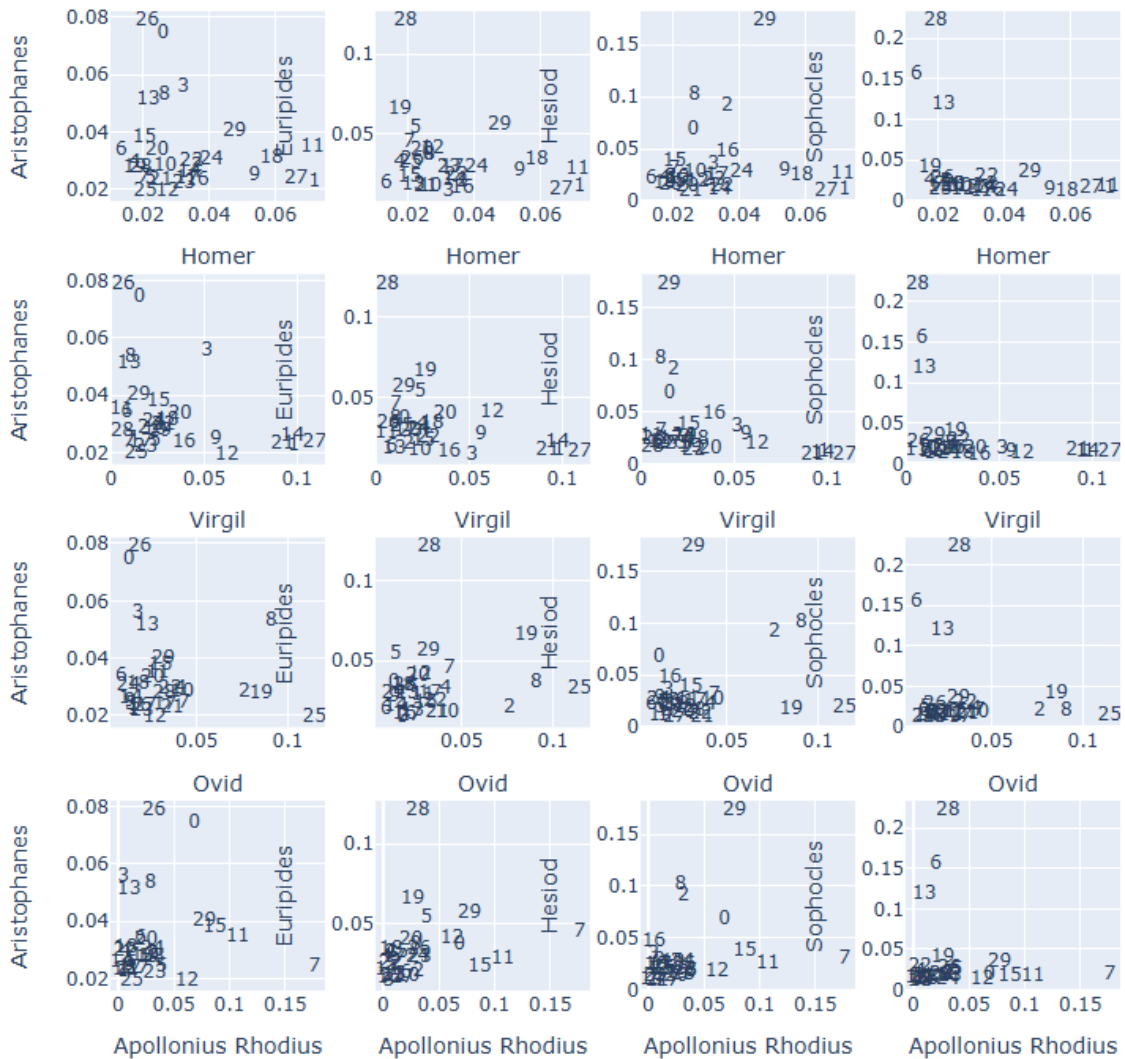


Figure 8: Word embedding pairwise comparison Epic vs Dramatic Poets.

With the word embedding relationships established, analysis of emotional patterns across the corpus through sentiment analysis provides additional context for understanding emotional ties to the vocabulary usage within each category. First, the VADER sentiment analysis of Homer's *Iliad* from the Epic category in Figure 9 shows the behavior of emotional progression throughout the work. The blue line represents the positive sentiment, and the orange line represents the negative sentiment. Overall, the positive sentiment dominates the beginning and some of the middle of the corpus, then there is a stark decline in the positive sentiment towards the end of the story. This behavior makes sense as Epic tales most often have tragic endings to round out the hero's journey, in this case with the brutal and bloody death of the strongest Trojan hero, Hector.

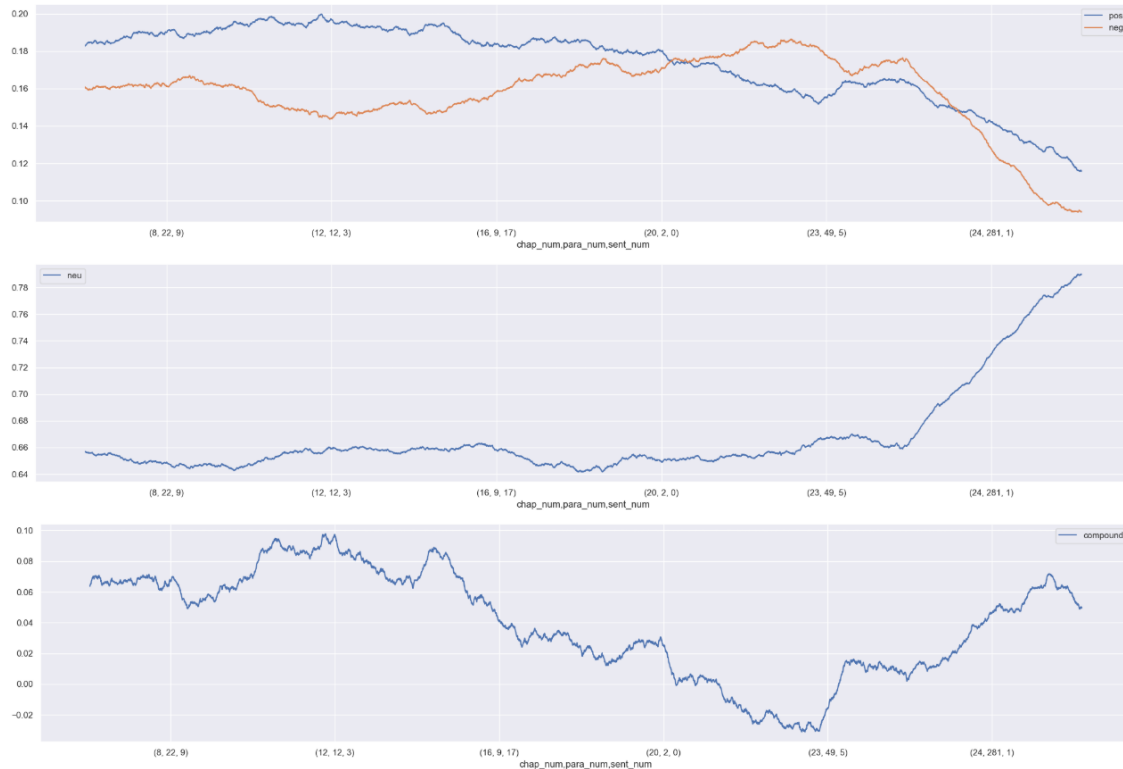


Figure 9: VADER sentiment analysis of Homer's Iliad.

Figure 10 is the Vader sentiment analysis of Euripides's *Medea*, part of the Dramatic Poetry category. The behavior of the positive and negative sentiment lines contrasts with the consistent progression of that in Figure 9. The positive sentiment in the story of *Medea* begins at a low point while the negative sentiment is at a high. From left to right there are two instances of positive and negative crossover and throughout the middle arch of the work the emotional theme crosses into positivity. Finally, the story closes with the second crossover back into negative sentiment. This behavior is consistent with Dramatic works as they start in dreary, bleak situations that the characters must preserve through. As the characters seemingly succeed in beating out the odds, there is a stark turn that befalls them, and the story ends in tragedy. The large separation in crossover points from negative to positive and the initial sentiment that starts the works are the key points of contrast between Epic and Dramatic Poetry works.

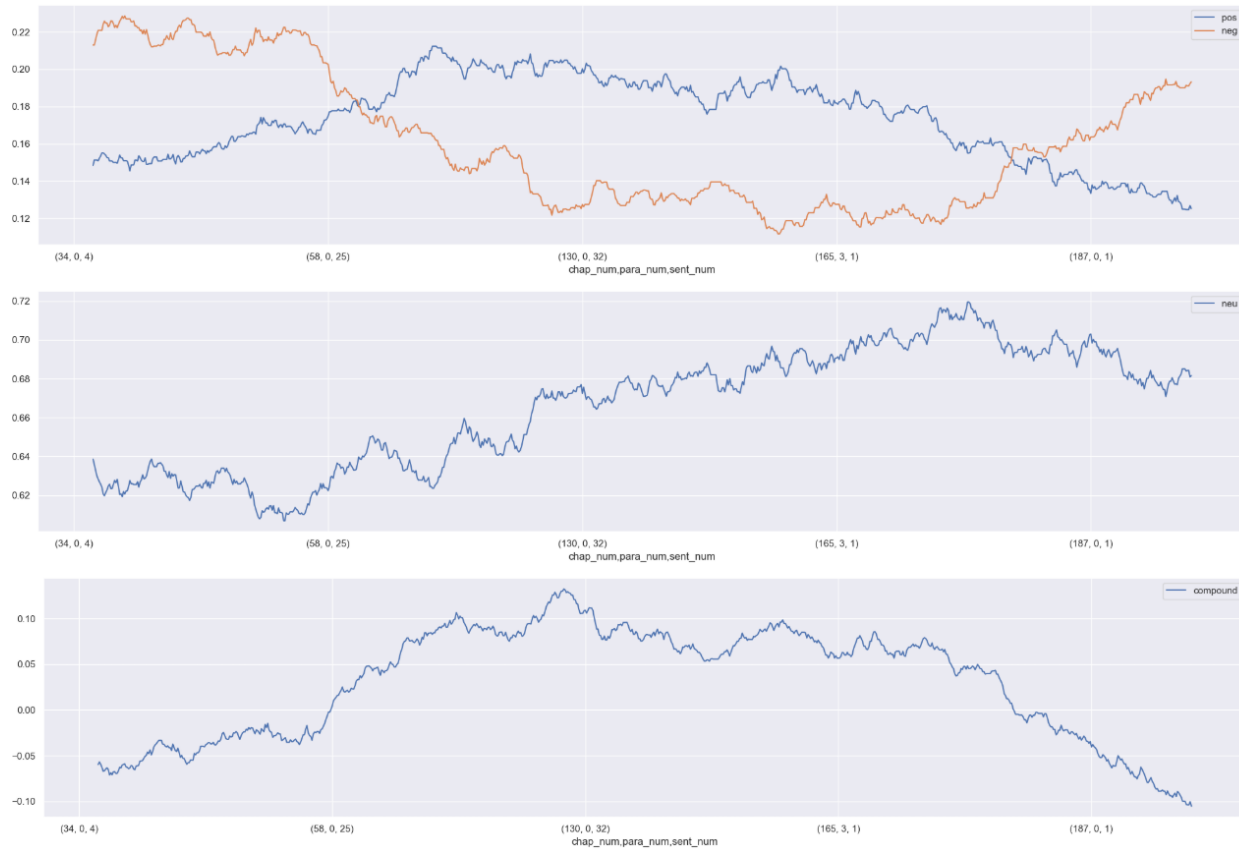


Figure 10: VADER sentiment analysis of Euripides' *Medea*.

Figure 11 is the Vader sentiment analysis of Hesiod's *Theogony*. The inclusion of this graph is to round out the category comparisons, as the emotional patterns of historical based texts are quite different than those of the classical Greco-Roman Epics or Dramatic Poetry. The vertical points of separation between the positive and negative sentiment lines are quite small in comparison to that of Figures 9 and 10. Additionally, the behavior of the sentiment lines here is much smoother when compared with the Figures 9 and 10, whose lines appear jagged from the incremental changes of highs and lows from left to right. This is likely because the mode of text expression is that of a historical narrative and not a dramatic one.

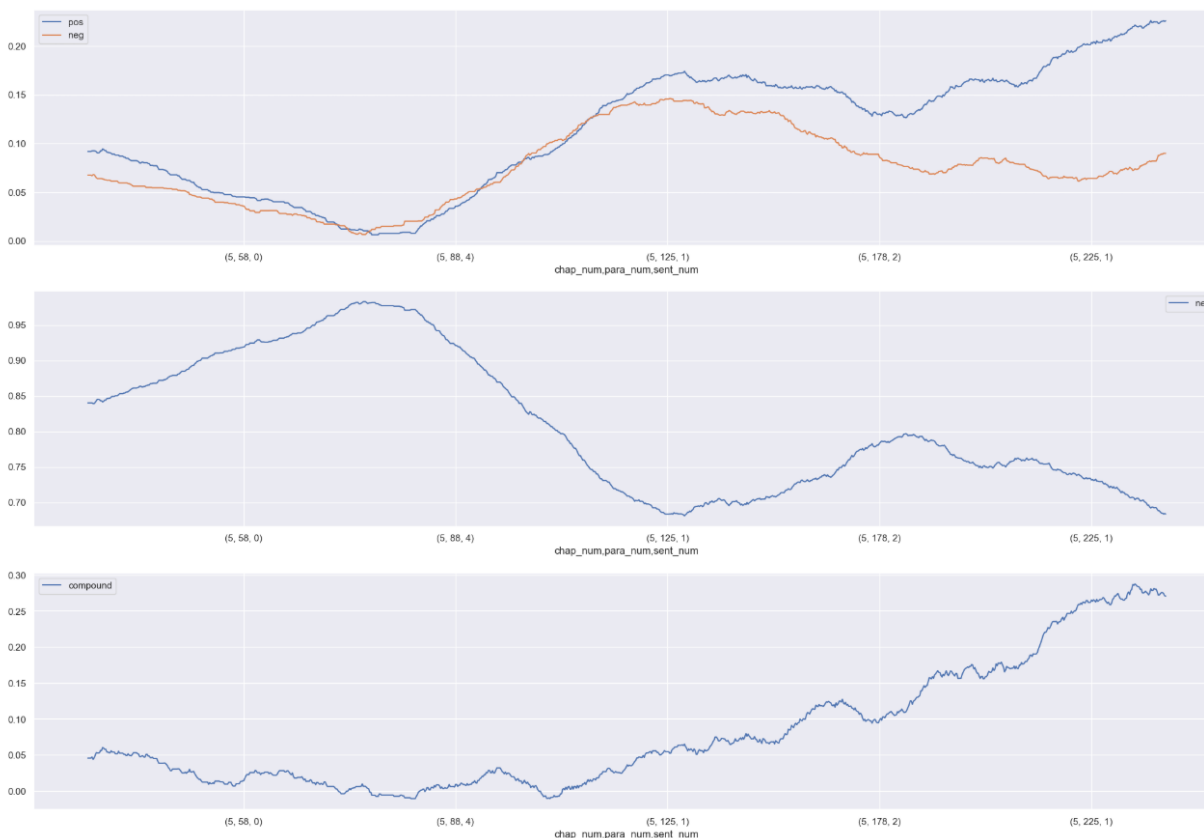


Figure 11: VADER sentiment analysis of the Hesiod's Theogony.

Overall, the contrasting patterns in the VADER graphs of positive and negative sentiment lines in the sentiment analysis provide quantitative evidence for distinguishing works into categories such as Epic, Dramatic Poetry, and History. These patterns highlight how distinct narrative purposes influence the shape and behavior of sentiment trajectories. For example, Epic works often align with themes of a hero's journey and feature steady emotional arcs. Dramatic Poetry exhibits sharp fluctuations reflecting rapid change in plot developments to evoke tension, and Historical works maintain consistent tones. These variations demonstrate how emotive patterns differ across categories.

Conclusion:

The analysis of the nine classical Greco-Roman corpus demonstrated patterns across the classic works categories of Epic, Dramatic Poetry, and Historical texts throughout varied ETA methods of analysis. The findings from the TFIDF analysis revealed the initial thematic/genre division of the sets of works in the corpus into the stated categories above. Where Epic works were strongly related to the vocabulary frequency of terms such as “gods”, “fate”, and “heroes”, Dramatic Poetry had terms like “chorus”, “revenge”, and “women” and historical text had more generic term usages. The thematic divisions were further supported with the LDA topic modeling and PCA analysis, which displayed distinct genre clustering and notable overlaps amongst

categories. For example, Apollonius's "The Argonautica," had genre elements of both the Epic and Dramatic Poetry category, which demonstrated that certain works could encompass elements of multiple genres. Lastly, the hierarchical clustering analysis brought another angle of support to the thematic division as there appeared to be three distinct clusters. The three clusters were comprised of vocabulary related to domestic life, interactions with the gods, and warfare, ultimately showing how these 3 clusters contain elements that are present across the corpus.

Furthermore, the sentiment analysis of the works revealed distinct emotive patterns related to each category/thematic genre. Epic works showed a behavior of an initial positive sentiment progression followed by a decline, alluding to the tragic ends of Epic works. Dramatic works, like Euripides Medea, exhibited sporadic changes in sentiment behavior, jagged lines, and several crossovers from positive to negative and negative to positive sentiments. Historical texts maintained a consistent emotional tone, as seen in the line behavior of the positive and negative sentiment, which were near smooth and consistent.

Finally, examining all the various analyses of the corpus gives meaningful evidence to support that classical works implement distinct narrative story telling strategies and emotional patterns that explore classical themes. These themes include interactions between mortals and gods, the hero's journey, and characters befalling tragedy. The analyses of these thematic patterns further suggest that classical works have distinct literary device applications that separate classical genres from one another. Further research with a larger corpus could reveal even more nuances to the patterns exhibited in classical works.

Deliverables Reference Guide:

All work is in a zip file called: final_project

All Deliverables are in a folder called: Project_Deliverables, but all work is in the above file location for ALL output files and tables

1. Collection of Source files called and located in: Source_Files located in Project_Deliverables
2. A manifest File Describing the source files: The top of the Jupyter Notebook contains the manifest information, heading called: Manifest Note
3. Collection of Data files called and located in: Data_Files located in Project_Deliverables
4. PCA Tables outputs called and located in: PCA_Outputs located in Project_Deliverables
5. LDA, word2vec, Sentiment Analysis: Located in the Final_Project.ipynb
6. Jupyter notebook that explores:
 - a. Hierarchical cluster diagrams
 - b. Heatmaps Showing correlations
 - c. Scatterplots
 - d. KDA plots
 - e. Dispersion plots
 - f. T-SNE plots
 - i. All are in the Final Project Jupyter Notebook called Final_Project.ipynb with write up of my findings.
7. Two-to-Four-page Final Report, this document.

Appendix

ⁱ <https://www.gutenberg.org/>