# Semantic flow in language networks

Edilson A. Corrêa Jr., Vanessa Q. Marinho, Diego R. Amancio

*Institute of Mathematics and Computer Science*
*University of São Paulo (USP)*
*São Carlos, São Paulo, Brazil*

## Abstract

In this study we propose a framework to characterize documents based on their semantic flow. The proposed framework encompasses a network-based model that connected sentences based on their semantic similarity. Semantic fields are detected using standard community detection methods. as the story unfolds, transitions between semantic fields are represent in Markov networks, which in turned are characterized via network motifs (subgraphs). Here we show that the proposed framework can be used to classify books according to their style and publication dates. Remarkably, even without a systematic optimization of parameters, philosophy and investigative books were discriminated with an accuracy rate of 92.5%. Because this model captures semantic features of texts, it could be used as an additional feature in traditional network-based models of texts that capture only syntactical/stylistic information, as it is the case of word adjacency (co-occurrence) networks.

*Keywords:* Complex Networks, Word Embeddings, Semantic Network, Text Similarity, Community Detection, Network Motifs

## 1. Introduction

In the last few years, several interesting findings have been reported by studies using network science to model language [1, 2, 3, 4]. Network-based models have been used e.g. to address the authorship recognition problem, where the structure of the networks

---

*Corresponding author
Email address:* diego@icmc.usp.br (Diego R. Amancio)

can provide valuable language-independent features. Other relevant applications relying on network science include the word sense disambiguation task [5, 6, 7], the analysis of text veracity and complexity [8]; and scientometric studies [9].

Whilst most of the network-based language research have been carried out at the word level [10, 11], only a limited amount of studies have been performed based on mesoscopic structures (sentences or paragraphs) [12]. In addition, most of the studies have analyzed language networks in a static way [13, 14]. In other words, once they are obtained, the order in which nodes (words, sentences, paragraphs) appear is disregarded. Here we probe the efficiency of sentence-based language networks in particular classification problems. Most importantly, differently from previous works hinging on network structure characterization [10, 11], we investigate whether the semantic flow along the narrative is an important feature for textual characterization in the considered classification tasks.

During the construction of a textual narrative, oftentimes authors follow a structured flow of ideas (introduction, narrative unfolding and conclusion). Even in books displaying a non-linear, complex narrative unfolding, one expects that an underlying linear semantic flow exists in authors' mind. In other words, even though narrative events might not organize themselves in a trivial linear form, the linearity imposed by written texts requires some type of linearization of the network of ideas. This idea is illustrated in Figure 1.

The ideas conveyed by a text can be represented as a complex network, where nodes represent semantic blocks (e.g. sentences, paragraphs), and edges are established according to semantic similarities. To map such a conceptual network into a text, authors perform a linearization process, where nodes (concepts, ideas) are linearly chosen and then transformed into a linear narrative (see Figure 1). Such a projection of a multidimensional space of ideas into a linear representation has been object of studies both on network theory and language research. A consequence of such a linearization in texts is the presence of long-range correlations at several linguistic levels, a property that has been extensively explored along the last years [16, 17, 18, 19].

While complex semantic networks have been used in previous works to represent the relationship between ideas and concepts, only a minor interest has been devoted to the
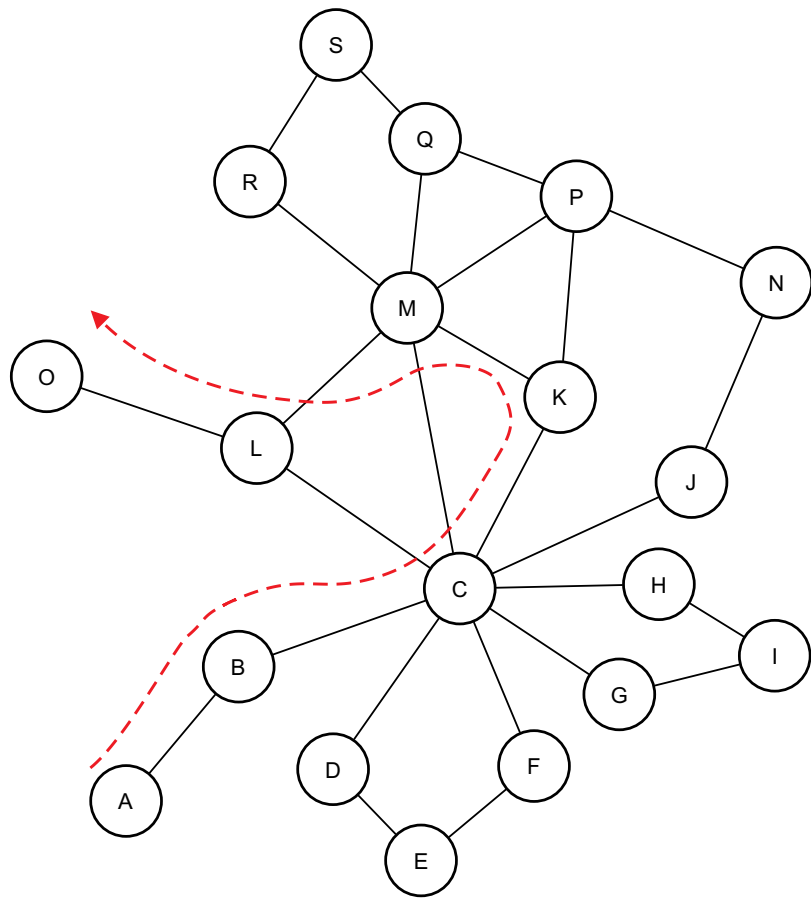
Figure 1: Illustration of process of construction a text from a network of a ideas. Each node represents an idea in this example and edges represent the relationship (similarity) of ideas. A written text can be seen as a walk on this network (see e.g. [15]). In the example, the following ideas are produced in the text: "A, B, C, K, M, L, 0".

analysis of how authors navigate the high-dimensional semantic relationships to generate a linear stream of words, sentences or paragraphs. In [20], a mesoscopic representation of networks was proposed. The authors used as a semantic, meaningful block a set of consecutive paragraphs. The semantic blocks were connected according to a lexical similarity index. The model aimed at combining a networked representation with a idea of semantic sequence obtained when reading a document. Even though some interesting patterns were found, the concept of semantic fields were not clear, as no semantic community structure arises from mesoscopic networks. The problem of linearization of a network structure was studied in [15]. A systematic analysis of the efficiency of several random walks in different topologies was probed. The efficiency was probed in a twofold manner: (i) the efficiency in transmitting the projected network; and (ii) the efficiency in recovering the original network. In [21], the authors explored the efficiency of navigating a idea space, by varying network topologies and exploration strategies.

In the current paper, we take the view that authors write documents by applying a linearization process to the original network of ideas, as shown in the procedure illustrated in Figure 1. Upon analyzing the flow of ideas with the adopted network-based framework, we show that features extracted from the networks can be employed to characterize and classify texts. More specifically, we defined the network of ideas as a network of sentences linked by semantic similarity. *Semantic fields* of similar sentences (nodes) were identified via network community detection. These fields (network communities) were then used to characterize the dynamics of authors' choices in moving from field to field as the story unfolds. Using a stochastic Markov model to represent the dynamics of choices of semantic fields performed by the author along the text, we showed, as a proof of principle, that the adopted representation can retrieve textual features including style (publication epoch) and complexity.

## 2. Research Questions

The main objective is to answer the following research questions: is there any patterns of semantic flow in stories? Are these patterns related to textual characteristics? To address

these questions, we use sentence networks to represent the semantic flow of ideas in texts. Such networks are summarized using a high-level representation based on the relationship between communities extracted from the sentence networks. Using this representation, we show that motifs extracted from such a high-level representation can be used to classify texts according to the style in which authors unfolds their stories. We argue that the obtained results suggest that the proposed high-level view of a text network could be further probed in other Natural Language Processing classification tasks.

This paper is organized as follows. Section 3 presents some concepts used in the proposed methods along with the method/framework itself. Section 4 presents the details of the experiments and results with a thorough discussion. Finally, in Section 5 we present some perspectives and insights for further works.

## 3. Materials and Methods

This study can be divided in two parts. In the first step, we identify the semantic clusters (fields) of the story. Differently from the analysis of short texts, where semantic groups can be identified mostly by identifying paragraphs, in long texts – the focus of this study – the identification of semantic clusters is more challenging because semantic topics might not be organized in consecutive sentences/paragraphs owing to the linearization process blueillustrated in Figure 1. In other words, the process of obtaining semantic clusters can be understood as the reverse operation depicted in Figure 1.

In order to identify semantic clusters from the text, we first create a network of sentences for each document, where sentences are linked if the similarity between them is above a given threshold. The obtained network is then analyzed via community detection methods, where groups of densely connected sentences are identified and considered as semantic clusters. A qualitative analysis of the obtained communities suggested that most of the largest communities are in fact related to a specific subtopic approached in the text (see details in Section 3.3). This idea relating semantic fields and network communities has also been used to construct automatic summarization systems [22].

In the second step of this study, we investigate the semantic flow of ideas developed by authors while unfolding their stories. We consider each community found as a semantic cluster, and as the story unfolds (one sentence after another), we analyze the community labels of the adjacent sentences to create a Markov chain, where each state represents a community and transitions are given by the text dynamics. Once the Markov chain representing the transitions of semantic clusters is obtained, the text is characterized by finding and counting different chain motifs. Such a characterization is then used to classify texts according to the semantic flow as revealed by sentences membership to different network communities.

The main objective of this work is to provide a framework to analyze and verify whether the semantic flow in texts can be used to characterize documents. Because the framework encompasses some steps, several alternatives could be probed in each step. We decided not to conduct a systematic analysis of combination of methods (and parameters) owing to the complexity of such analysis. A systematic study of the parameters and methods optimizing the proposed framework is intended to be conducted as a future work.

In Figure 2 we show a representation of the framework proposed to analyze stories. In the next section, we detail each of steps used in this framework.

*3.1. Word and Sentence Embeddings*

Usually any vector representation of words is known as a *word embedding*. However, since the creation of *neural word embeddings* [23], the term is mostly used to name those approaches based on neural network representations. The *word embedding* model proposed in [23] aimed at classifying texts based on raw text input. Thus, the classification does not require that textual features as input. Typically, *word embeddings* are dense vectors that are learned for a specific vocabulary, with the objective of addressing some task.

A typical task addressed with word embeddings is the language modeling problem, which aims at learning a probability function describing the sequence of words in a language. More recently, this same vector representation has been used in more complex models, with the objective of addressing several Natural Language Processing tasks simultaneously, including POS tagging, name entity recognition, semantic role labeling and others [24, 25]. Despite
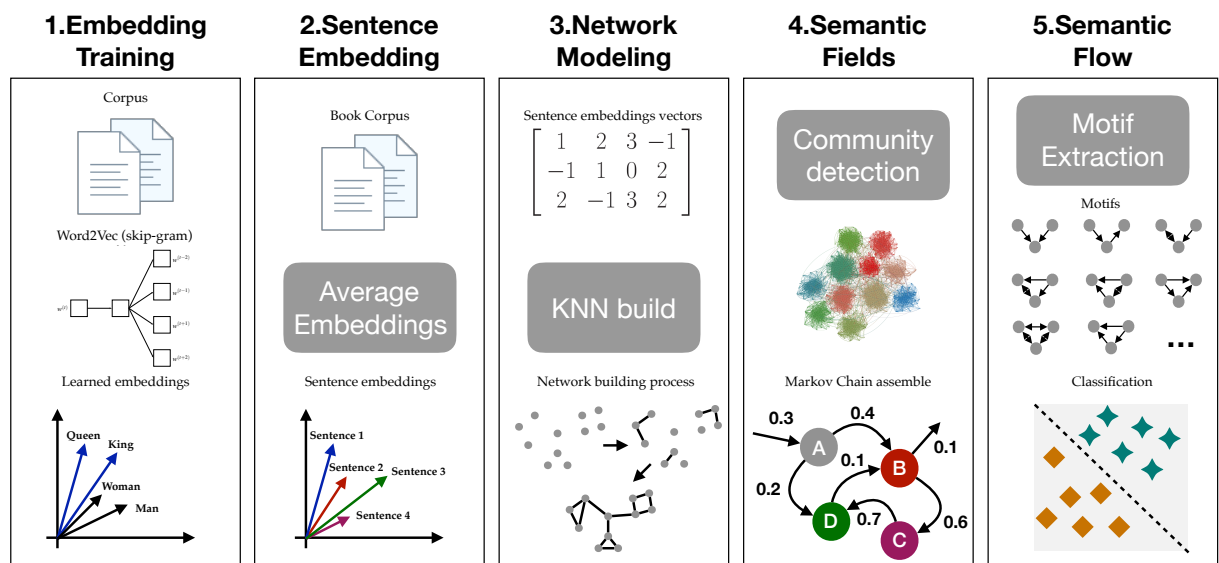
6

Figure 2: Sequence of steps employed to characterize documents using the proposed framework: (1) word embedding generation; (2) sentence embeddings generation from word embeddings; (3) a sentence similarity network is generated based on the similarity of sentence embeddings; (4) network communities are detected and a Markov chain is built based on the story unfolding (semantic flow); and (5) motifs are identified in the Markov chain representing the semantic flow. These motifs are then used as features in a classification method.

its relative success in the above mentioned tasks, the adopted embeddings could not be used in general purpose applications [24, 25]. In order to allow the use of embeddings in wider contexts, the Word2Vec representation was proposed [26, 27].

The Word2Vec is a neural model proposed to learn a dense, high-quality representation that is able to capture both syntactical and semantical language properties. As a consequence, vectors representing words conveying the same meaning are close in the considered space. An interesting property of the Word2Vec technique is the *compositionality*, which allows that large information blocks (e.g. sentences) can be represented by combining the representation of the vector representing the words in the sentence. Other interesting property is the ability to combine embeddings in a intuitive fashion [26, 28]. For example, using the Word2Vec technique, the following relationship can be obtained:

$$\text{vector(``King'') - vector(``Man'') + vector(``Woman'')} \simeq \text{vector(``Queen'').} \tag{1}$$

The Word2Vec model is a robust, general-purpose neural representation that has been widely used in several Natural Language Processing tasks, including machine translation [29], summarization [30, 31], sentiment analysis [32] and others. Given the success of the this model and the possibility of composition in different scenarios (sentiment analysis and sense disambiguation) [32, 33, 34], in the current study we used a representation of sentences based on the Word2Vec. More specifically, here the embedding $\mathbf{s}$ of a sentence $s$ is represented by the average embedding of the words in $s$:

$$\mathbf{s} = \frac{1}{\omega(s)} \sum_{i=0}^{\omega} \mathbf{w}_i. \tag{2}$$

where $\mathbf{w}_i$ is the embedding of the $i$-th word in $s$ and $\omega(s)$ is the total number of words in $s$.

The word embedding technique used here was obtained with the Word2Vec method (skip-gram). The training phase used the Google News corpus [27, 26]. According to [27, 26], the parameters of the method are optimized in the context of semantical similarity task. The combination of embeddings to represent a sentence in equation 2 could also be performed by summing individual embeddings. However, it has been shown that there is no significant difference when sentence embeddings are used to construct a network of sentence

similarity [35]. We note that some words are removed from this analysis. This includes *stopwords* (e.g. articles and prepositions) and words with no embeddings in the Google News corpus. Thus, whenever a sentence contains only words with no available embeddings, it is removed from the analysis.

### 3.2. Modeling sentence embeddings into complex networks

This step corresponds to the reverse process illustrated in Figure 1. In other words, a network representing the relationship between ideas is created from the text. The construction of networks from vector structures has been explored in recent works. In [36], the authors present such a transformation as a framework in complex systems analysis. The transformation of vector structures into networks has also been used in the context of text analysis [35, 37]. The creation of a complex network from Word2Vec was proposed using a twofold approach. The $d$-proximity technique links all nodes whose distance from the refence node is lower than $d$. The second technique is the $k$-NN approach, which links all $k$ nearest nodes to the reference node. In the same line, [35] created a network based on word embeddings. However, the authors aimed at creating a network that takes into account the sense of words to solve ambiguities. Each occurrence of an ambiguous words was modelled as a node in the network. Nodes were represented by a vector combining the embeddings of the words in the context. Two occurrences of an ambiguous were then connected whenever the respective embeddings were similar. In other words, two ambiguous words were linked if they appeared in similar contexts.

In the currrent study, sentences were connected according to the $k$-NN technique, as suggested by other works [37]. Each sentence is represented as a vector according to equation 2. The value of $k$ in the main experiments were chosen to allow that each network is composed of a single connected component. In particular, the lowest $k$ allowing the creation of a connected network was used for each book.

### 3.3. Community detection

The next step in the proposed framework concerns the detection of semantic fields, i.e. the communities in the network of sentences. A recurrent phenomena in several complex

networks is the existence of communities, i.e. groups of strongly connected nodes. Similarly to other network measurements, the detection of communities gives important information regarding the organization of networks. Communities are present in different networks including in biological, social and information networks [38].

A well-known measure to quantify the quality of partitions in complex networks is the modularity [39, 40]. This measure compares the obtained partition with a null model, i.e. a network with similar properties but with no community structure. Several algorithms have been proposed to address the community detection problem via optimization of the modularity. In the main experiments we used the the Louvain method [41] to identify communities. The main advantage of this method is its computational efficiency, which has allowed its use in several contexts [37, 42]. Another advantage associated to this algorithm is that no additional parameters are required to optimize the modularity.

In the proposed network representation, communities represent groups of interconnected sentences about a given topic. Because the $k$-NN construction allows nodes to be connected to other close nodes and, considering the Word2Vec an efficient semantic representation, the linking strategy allows the creation of dense clusters of semantically related sentences. This idea of semantic clusters has also been explored via community detection in similar works [22, 35, 37]. For example, using networks built at the word level, the groups detected in [37] were found to represent large cities, professions and others topics. In [35], the obtained groups were found to represent words conveying the same sense.

In order to illustrate the process of obtaining semantic communities, we performed an analysis of the obtained communities in the book "Alice's Adventures in Wonderland", by Lewis Carroll. We summarize below the main topics approaches in some of the communities obtained by the Louvain algorithm:

1. *Community A*: this community includes sentences mentioning animals (e.g. "pet", "cat", "mouse" and "dog"). This community also includes dialogues between Alice and animals. "Cat" is the main character in this community.

2. *Community B*: this community includes words sentiment words expressed via speeches.

Some of the words in this community are "passionate", "melancholy", "angrily", "shouted" and "screamed".

3. *Community C*: this community includes several adverbs related to Alice's actions.

4. *Community D*: this community includes words related to sentiments such as anger, tranquility and peacefulness.

5. *Community E*: this community is most related to the word "soup".

6. *Community F*: this community is related to geographical locations, including countries and cities (Australia, Rome and New Zealand). Interestingly, this community also included the word "Cricket", a prominent sport in Australia.

7. *Community G*: this community included mostly sentences referring to "Dormouse", one of the main characters in the plot.

While most of the obtained communities are informative, a few communities were found to be more dispersed, approaching more than one topic. This might occur given the limitations of the embeddings model, since some words might not be available in the considered model. Despite these limitations, we show that the flow of information (from sentence to sentence) in the obtained semantic communities can be used to characterize texts.

*3.4. Markov Chains*

In order to capture how authors move from community to community (semantic field) as their story unfolds, we create a representation of community transitions. The idea of studying language via Markov process is not recent. One of the first uses of this model is the study of letters sequences [43]. Since then, Markov chains are used as a statistical tool in several natural language processing problems, including language modeling, machine translation and speech recognition [44].

Here we represented the transitions between semantic fields (network communities) as a first order Markov chain. In this representation, each community becomes a state. Note that this approach of representing communities as a single unit has also been used in other contexts [9]. The probabilities of transition are considered according to the frequency of transitions observed in adjacent sentences. As we shall show, using this model, it is possible

Figure 3: Example of sentence network obtained from the book "Alice's Adventures in Wonderland", by Lewis Carroll. Colors represent community labels obtained with the Louvain method. The visualization was obtained with the method described in [9].

to detect patterns (see Section 3.5) of how authors change topics in their stories. As a proof of principle, these patterns are used to characterize texts in distinct classification tasks.

The process of creating a Markov chain from a network divided into communities is shown in Figure 4. In the previous phase, communities are identified to represent distinct semantic field of the story (see left graph in Figure 4). Because each sentence belongs to just one community, the text can be regarded as discrete time series, where each element corresponds to the membership (community label) of each sentence. Using this sequence of community labels, it is possible to create a Markov chain representing all transitions between communities (see graph on the top left of Figure 4). Transitions weights are proportional to the frequency in which they occur and normalized so as to represent a probability. This representation is akin to a Markov chain used in other works addressing the language modeling problem [45]. The main difference here is that we are not interested in the use of particular words, but in semantic fields [46]. Once the Markov chain is obtained, we characterize this structure using *network motifs* (see Section 3.5).

## 3.5. Motifs

Network motifs are used to analyze a wide range of complex systems, including in biological, social and information networks [47]. Motifs can be defined as small subgraphs occurring in real systems in a significant way. To quantify the significance, in general, one assumes an equivalent random network as null model. In text analysis, motifs have been used to analyze word networks [48] in applications focusing on the syntax and style of texts. More recently, an approach based on labelled motifs showed that authors tend to use words in combination with particular motifs [49].

While the structure of the Markov Chains could be analyzed using traditional network measurements, we decided not to use these measurements owing to the limited size of these structures. As suggested in related works, a characterization based on network metrics in small networks might not be informative [50, 14, 3]. As we shall show in the results, this is a simple, yet useful approach to classify small Markov Chains.

Three different approaches were considered to extract motifs from Markov networks.

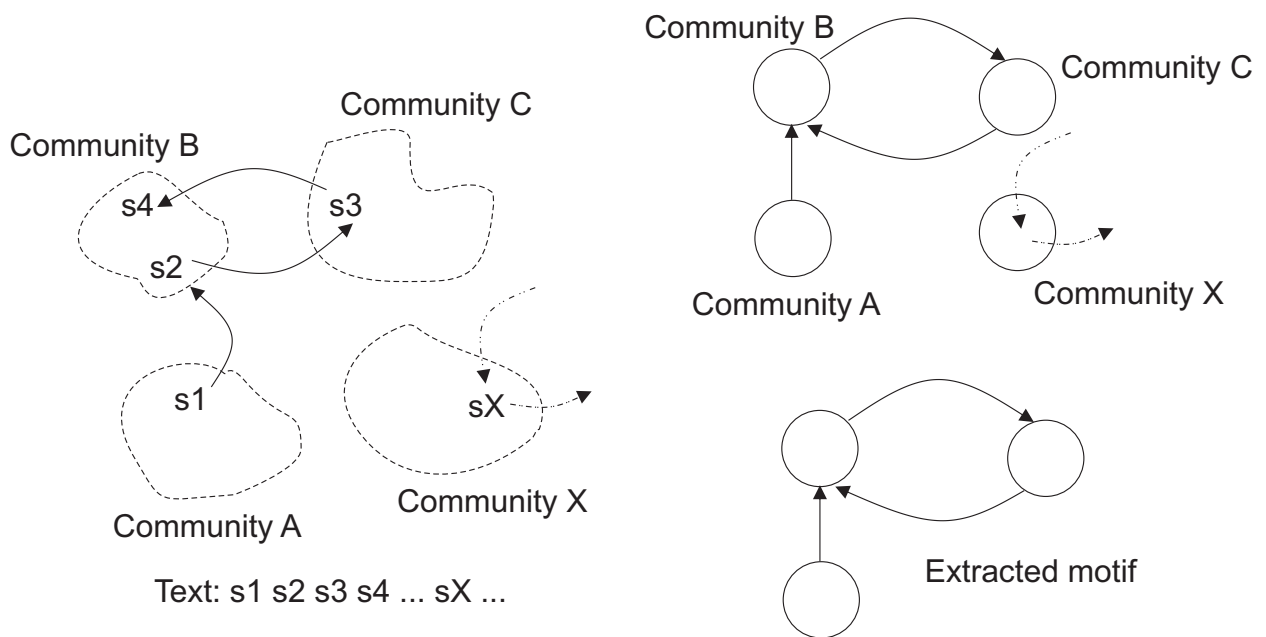Figure 4: Example of extraction of motifs from the network. As the text unfolds according to a given order of sentences $(s1, s2, s3, s4 \ldots sX \ldots)$ a sequence of communities is generated (Community A, Community B, Community C, Community B). This sequence is used to create a Markov chain. Finally, the Markov chain is characterized by counting different patterns (motifs) of community transitions.

We discriminated strategies according to the use of weights to count motifs (unweighted vs weighted). If a thresholding is applied before extracting motifs, the strategy is named with the "simplified":

1. *Unweighted strategy*: no thresholding is applied. All weights are disregarded. Every time a motif is detected, its frequency is increased by one.

2. *Simplified unweighted strategy*: this approach is the same as the *unweighted* strategy. However, before counting motifs, the weakest edges are removed according to a given threshold.

3. *Simplified weighted strategy*: before counting motifs, the weakest edges are removed according to a given threshold. All edges weights are considered in the remaining Markov network. Every time a given motif is found, the respective "frequency" of that motif is increased by the sum of the weights of its edges.

### 3.6. Classification - Machine Learning Methods

The extracted motifs from the Markov Chains are used as input (features) to the classification systems. The following methods were used in the experiments: Decision Tree (CART), kNN, SVM (linear) and Naive Bayes [51]. The evaluation was performed using a 10-fold cross-validation approach. As suggested in related works, all classifiers were trained with their default configuration of parameters [52].

## 4. Results and Discussion

### 4.1. Classification tasks

Here we probed whether the dynamics of changes in semantic groups in books can be used to characterize stories. The proposed methodology was applied in two distinct classification tasks. In the first task, we aimed at distinguishing three different thematic classes: (i) children books; (ii) investigative; and (iii) philosophy books. The second aimed at discriminating books according to their publication dates. All books (and their respective classes) were obtained from the Gutenberg repository. The list of books and respective authors are listed in the Supplementary Information.

In the first experiment, we evaluated if patterns of semantic changes are able to distinguish between children, philosophy or investigative books. We considered problems with two or three classes. The obtained results are shown in Table 1. In this case, weights were disregarded after the construction of the Markov networks (*unweighted* version). Considering subtasks encompassing only two classes, only the distinction between children and investigative texts were not significant, with a low accuracy rate. The distinction philosophy books and the other two classes, however, yielded a much better discrimination. These results were found to be significant. When all three classes are discriminated, a low accuracy rate was found (43.8%), even though this still represents a significant result. The low accuracy rate found using the proposed approach is a consequence of a regular behavior found in the Markov chains. In other words, in most of the books, all communities were found to be connected to each other, hampering thus the discriminability of different types of books.

Table 1: Accuracy rate and $p$-value obtained for the classification substaks. Only the best results are shown among all considered classifiers. We considered the *unweighted* version of the Markov networks to extract motifs.

| Subtask | Acc. | $p$-value |
|---|---|---|
| children $\times$ investigative | 50.8% | $5.56 \times 10^{-1}$ |
| children $\times$ philosophy | 71.6% | $1.30 \times 10^{-3}$ |
| investigative $\times$ philosophy | 70.8% | $3.30 \times 10^{-3}$ |
| children $\times$ investigative $\times$ philosophy | 43.8% | $3.50 \times 10^{-2}$ |

Given the low accuracy rates obtained with the *unweighted* strategy, we analyzed if the *simplified* unweighted version was able to provide a better characterization. In this case, the weakest edges were removed before the extraction of motifs. We considered the thresholding ranging between 0.01 and 0.20. The main idea here is to remove less important links between communities. The obtained results are shown in Table 2. All obtained results turned out to be significant. All previous accuracy rates were improved. Interestingly, a high discrimination rate (91.6%) was obtained when discriminating investigative and philosophy books. These results suggest that the threshold is an important pre-processing step here,

given that it can boost the performance of the classification by a large margin.

Table 2: Accuracy rate and $p$-value obtained for the classification substaks. Only the best results are shown among all considered classifiers and thresholds. We considered the *simplified unweighted* version of the Markov networks to extract motifs.

| Subtask | Acc. | Threshold | $p$-value |
|---|---|---|---|
| children $\times$ investigative | 65.8% | 0.060 | $1.64 \times 10^{-2}$ |
| children $\times$ philosophy | 81.0% | 0.190 | $1.19 \times 10^{-5}$ |
| investigative $\times$ philosophy | 91.6% | 0.075 | $2.23 \times 10^{-10}$ |
| children $\times$ investigative $\times$ philosophy | 62.2% | 0.075 | $2.00 \times 10^{-7}$ |

When combining thresholding and edges weights in the *simplified weighted* version, the results obtained in Table 3 were further improved. The highest gain in performance was observed when discriminating children from philosophy books: the performance improved from 81.0% to 89.0%. Only a minor improvement was observed when all three classes were discriminated. Overall, this results suggest that both thresholding and the use of edges weights might be useful to characterize Markov networks. Most importantly, all three methods showed that, in fact, there is a correlation between the thematic approached and the way in which authors approaches semantic groups in texts.

Table 3: Accuracy rate and $p$-value obtained for the classification substaks. Only the best results are shown among all considered classifiers and thresholds. We considered the *simplified weighted* version of the Markov networks to extract motifs.

| Subtask | Acc. | Threshold | $p$-value |
|---|---|---|---|
| children $\times$ investigative | 70.8% | 0.075 | $3.30 \times 10^{-3}$ |
| children $\times$ philosophy | 89.0% | 0.145 | $1.62 \times 10^{-8}$ |
| investigative $\times$ philosophy | 92.5% | 0.120 | $2.23 \times 10^{-10}$ |
| children $\times$ investigative $\times$ philosophy | 62.7% | 0.075 | $2.00 \times 10^{-7}$ |

We also investigated if the patterns of semantic flow varies with the publication date. For this reason, we selected a dataset with books in different periods. The following classes were considered, according to the range of publication dates:

17

1. Books published between 1700 and 1799 (30 books in each class).

2. Books published between 1800 and 1899 (30 books in each class).

3. Books published after 1900 (30 books in each class).

4. Books published between 1700 and 1850 (45 books in each class).

5. Books published after 1851 (45 books in each class).

The results obtained in the classification for different subtasks is shown in Table 4. We only show here the results obtained for the simplified unweighted characterization because it yielded the best results. Overall, all classification results are significant, confirming thus that there are statistically significant differences of semantic flow patterns for books published in different epochs. However, the results obtained here are worse than the ones obtained in the dataset with books about different themes (see Table 3). Therefore, patterns of semantic flow seems to be less affected by the year of publication, while being more sensitive to the subject/topic approached by the text.

Table 4: Performance of the proposed method using the *simplified unweighted* motif characterization of Markov networks. For each subtask, only the best threshold obtained for the best classifier is shown.

| Subtask | Acc. | Threshold | $p$-value |
|---|---|---|---|
| $1700 - 1799 \times 1800 - 1899$ | 70.0% | 0.195 | $1.34 \times 10^{-3}$ |
| $1700 - 1799 \times 1900$ or later | 75.0% | 0.060 | $6.73 \times 10^{-5}$ |
| $1800 - 1899 \times 1900$ or later | 70.0% | 0.160 | $1.34 \times 10^{-3}$ |
| $1700 - 1850 \times 1851$ or later | 66.0% | 0.010 | $6.74 \times 10^{-3}$ |
| $1700 - 1799 \times 1800 - 1899 \times 1900$ or later | 55.0% | 0.025 | $1.22 \times 10^{-5}$ |

## 5. Conclusion

In this paper we investigate whether patterns of semantic flow arises for different classes of texts. To represent the relationship between ideas in texts, we used a sentence network representation, where sentences (nodes) are connected based on their semantic similarity. Semantic clusters were identified via community detection and high-level representation of

each book was created based on the transition between communities as the story unfolds. Finally, motifs were extracted to characterize the patterns of transition between semantic groups (communities). When applied in two distinct tasks, interesting results were found. In the task aiming at classifying books according to the approached themes, we found an high accuracy rate (92.5%) when discriminating investigative and philosophy books. A significant performance in the classification was also obtained when discriminating books published in distinct epochs. However, the discriminability for this task was not as high as the ones obtained when discriminating investigative, philosophy and children books.

Given the complexity of the components in the proposed framework, we decided not to optimize each step of the process. Even without a rigorous optimization process, we were able to identify semantic flow patterns that were able to discriminate distinct classes of texts. As future works, we intend to perform a systematic analysis on how to optimize the process. For example, during the construction of the networks, different approaches could be used to link similar sentences. In a similar fashion, different strategies to identify communities could also be used in the analysis. Finally, we could also investigate additional approaches to characterize the obtained Markov networks.

The results obtained here suggests that different classes of books display different patterns of semantic flow. This suggests that the semantic flow could play an important role in other NLP tasks. For example, in the authorship recognition task, patterns extracted from a semantic flow analysis could be combined with other techniques to improve the characterization of authors. A similar idea could also be applied to the analysis of other stylometric tasks. Since semantic networks have been studied in cognitive sciences, we believe that the adopted network representation could be adapted and used – as an auxiliary tool – to study complex brain and cognitive processes that could assist the diagnosis of cognitive disorders via text analysis [53, 54].

**Acknowledgements**

## References

## References

[1] W. Jin, R. K. Srihari, Graph-based text representation and knowledge discovery, in: Proceedings of the 2007 ACM symposium on Applied computing, ACM, 2007, pp. 807–811.

[2] R. F. Cancho, R. V. Solé, R. Köhler, Patterns in syntactic dependency networks, Physical Review E 69 (5) (2004) 051915.

[3] D. R. Amancio, Comparing the topological properties of real and artificially generated scientific manuscripts, Scientometrics 105 (3) (2015) 1763–1779.

[4] M. A. Montemurro, D. H. Zanette, Keywords and co-occurrence patterns in the voynich manuscript: An information-theoretic analysis, PloS one 8 (6) (2013) e66344.

[5] E. Agirre, A. Soroa, Personalizing pagerank for word sense disambiguation, in: Proceedings of the 12th Conference of the European Chapter of the Association for Computational Linguistics, Association for Computational Linguistics, 2009, pp. 33–41.

[6] T. C. Silva, D. R. Amancio, Word sense disambiguation via high order of learning in complex networks, EPL (Europhysics Letters) 98 (5) (2012) 58001.

[7] D. R. Amancio, O. N. Oliveira Jr., L. F. Costa, Unveiling the relationship between complex networks metrics and word senses, EPL (Europhysics Letters) 98 (1) (2012) 18002.

[8] D. R. Amancio, O. N. Oliveira Jr, L. da Fontoura Costa, Identification of literary movements using complex networks to represent texts, New Journal of Physics 14 (4) (2012) 043029.

[9] F. N. Silva, D. R. Amancio, M. Bardosova, L. d. F. Costa, O. N. Oliveira, Using network science and text analytics to produce surveys in a scientific topic, Journal of Informetrics 10 (2) (2016) 487–502.

[10] R. F. Cancho, R. V. Solé, The small world of human language, Proceedings of the Royal Society of London B: Biological Sciences 268 (1482) (2001) 2261–2265.

[11] H. Liu, J. Cong, Language clustering with word co-occurrence networks based on parallel texts, Chinese Science Bulletin 58 (10) (2013) 1139–1144.

[12] H. F. Arruda, V. Q. Marinho, L. F. Costa, D. R. Amancio, Paragraph-based representation of texts: a complex networks approach, Information Processing and Management 56 (2019) 479–494.

[13] C. Akimushkin, D. R. Amancio, O. N. Oliveira Jr, Text authorship identified using the dynamics of word co-occurrence networks, PloS one 12 (1) (2017) e0170527.

[14] D. R. Amancio, Probing the topological properties of complex networks modeling short written texts, PLoS ONE 10 (2) (2015) e0118394.

[15] H. F. de Arruda, F. N. Silva, C. H. Comin, D. R. Amancio, L. F. Costa, Connecting network science and information theory, Physica A: Statistical Mechanics and its Applications 515 (2019) 641–648.

[16] W. Ebeling, A. Neiman, Long-range correlations between letters and sentences in texts, Physica A: Statistical Mechanics and its Applications 215 (3) (1995) 233 – 241.

[17] A. Schenkel, J. Zhang, Y.-C. Zhang, Long range correlation in human writings, Fractals 1 (01) (1993) 47–57.

[18] E. Alvarez-Lacalle, B. Dorow, J.-P. Eckmann, E. Moses, Hierarchical structures induce long-range dynamical correlations in written texts, Proceedings of the National Academy of Sciences 103 (21) (2006) 7956–7961.

[19] M. Amit, Y. Shmerler, E. Eisenberg, M. Abraham, N. Shnerb, Language and codification dependence of long-range correlations in texts, Fractals 2 (01) (1994) 7–13.

[20] H. F. Arruda, F. N. Silva, V. Q. Marinho, D. R. Amancio, L. F. Costa, Representation of texts as complex networks: a mesoscopic approach, Journal of Complex Networks 6 (1) (2018) 125–144.

[21] H. F. Arruda, F. N. Silva, L. F. Costa, D. R. Amancio, Knowledge acquisition: A complex networks approach, Information Sciences 421 (2017) 154–166.

[22] L. Antiqueira, O. N. Oliveira Jr, L. da Fontoura Costa, M. d. G. V. Nunes, A complex network approach to text summarization, Information Sciences 179 (5) (2009) 584–599.

[23] Y. Bengio, R. Ducharme, P. Vincent, C. Jauvin, A neural probabilistic language model, journal of machine learning research 3 (Feb) (2003) 1137–1155.

[24] R. Collobert, J. Weston, A unified architecture for natural language processing: Deep neural networks with multitask learning, in: Proceedings of the 25th international conference on Machine learning, ACM, 2008, pp. 160–167.

[25] R. Collobert, J. Weston, L. Bottou, M. Karlen, K. Kavukcuoglu, P. Kuksa, Natural language processing (almost) from scratch, Journal of Machine Learning Research 12 (Aug) (2011) 2493–2537.

[26] T. Mikolov, K. Chen, G. Corrado, J. Dean, Efficient estimation of word representations in vector space, arXiv preprint arXiv:1301.3781.

[27] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, J. Dean, Distributed representations of words and phrases and their compositionality, in: Advances in neural information processing systems, 2013, pp. 3111–3119.

[28] T. Mikolov, W.-t. Yih, G. Zweig, Linguistic regularities in continuous space word representations, in: Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, 2013, pp. 746–751.

[29] I. Sutskever, O. Vinyals, Q. V. Le, Sequence to sequence learning with neural networks, in: Advances in neural information processing systems, 2014, pp. 3104–3112.

[30] A. M. Rush, S. Chopra, J. Weston, A neural attention model for abstractive sentence summarization, in: Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, 2015, pp. 379–389.

[31] J. V. Tohalino, D. R. Amancio, Extractive multi-document summarization using multilayer networks, Physica A: Statistical Mechanics and its Applications 503 (2018) 526 – 539.

[32] R. Socher, A. Perelygin, J. Wu, J. Chuang, C. D. Manning, A. Ng, C. Potts, Recursive deep models for semantic compositionality over a sentiment treebank, in: Proceedings of the 2013 conference on empirical methods in natural language processing, 2013, pp. 1631–1642.

[33] E. A. Corrêa Júnior, V. Q. Marinho, L. B. dos Santos, Nilc-usp at semeval-2017 task 4: A multi-view ensemble for twitter sentiment analysis, in: Proceedings of the 11th International Workshop on Semantic Evaluation (SemEval-2017), Association for Computational Linguistics, 2017, pp. 611–615.

[34] I. Iacobacci, M. T. Pilehvar, R. Navigli, Embeddings for word sense disambiguation: An evaluation study, in: Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), Association for Computational Linguistics, Berlin, Germany, 2016, pp. 897–907.

[35] E. A. Corrêa Jr, D. R. Amancio, Word sense induction using word embeddings and community detection in complex networks, arXiv preprint arXiv:1803.08476.

[36] C. H. Comin, T. K. Peron, F. N. Silva, D. R. Amancio, F. A. Rodrigues, L. d. F. Costa, Complex systems: features, similarity and connectivity, arXiv preprint arXiv:1606.05400.

[37] B. Perozzi, R. Al-Rfou, V. Kulkarni, S. Skiena, Inducing language networks from continuous space word representations, in: Complex Networks V, Springer, 2014, pp. 261–273.

[38] S. Fortunato, Community detection in graphs, Physics Reports 486 (3) (2010) 75 – 174.

[39] M. E. Newman, M. Girvan, Finding and evaluating community structure in networks, Physical review E 69 (2) (2004) 026113.

[40] M. E. Newman, Modularity and community structure in networks, Proceedings of the national academy of sciences 103 (23) (2006) 8577–8582.

[41] V. D. Blondel, J.-L. Guillaume, R. Lambiotte, E. Lefebvre, Fast unfolding of communities in large networks, Journal of statistical mechanics: theory and experiment 2008 (10) (2008) P10008.

[42] E. Corrêa Jr., A. A. Lopes, D. R. Amancio, Word sense disambiguation: A complex network approach, Information Sciences 442-443 (2018) 103 – 113.

[43] A. A. Markov, An example of statistical investigation of the text eugene onegin concerning the connection of samples in chains, Vol. 7, 1913, pp. 153–162.

[44] C. D. Manning, H. Schütze, Foundations of Statistical Natural Language Processing, MIT Press, Cambridge, MA, USA, 1999.

[45] J. M. Ponte, W. B. Croft, A language modeling approach to information retrieval, in: Proceedings of the 21st annual international ACM SIGIR conference on Research and development in information retrieval, ACM, 1998, pp. 275–281.

[46] F. Li, T. Dong, Text categorization based on semantic cluster-hidden markov models, in: International Conference in Swarm Intelligence, Springer, 2013, pp. 200–207.

[47] R. Milo, S. Shen-Orr, S. Itzkovitz, N. Kashtan, D. Chklovskii, U. Alon, Network motifs: simple building blocks of complex networks., Science 298 (5594) (2002) 824–827.

[48] D. R. Amancio, E. G. Altmann, D. Rybski, O. N. Oliveira Jr, L. F. Costa, Probing the statistical properties of unknown texts: application to the voynich manuscript, PLoS One 8 (7) (2013) e67310.

[49] V. Q. Marinho, G. Hirst, D. R. Amancio, Labelled network subgraphs reveal stylistic subtleties in written texts, Journal of Complex Networks 6 (4) (2018) 620–638.

[50] B. C. Van Wijk, C. J. Stam, A. Daffertshofer, Comparing brain networks of different size and connectivity density using graph theory, PloS one 5 (10) (2010) e13701.

[51] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, et al., Scikit-learn: Machine learning in python, Journal of machine learning research 12 (Oct) (2011) 2825–2830.

[52] D. R. Amancio, C. H. Comin, D. Casanova, G. Travieso, O. M. Bruno, F. A. Rodrigues, L. da Fontoura Costa, A systematic comparison of supervised classifiers, PloS one 9 (4) (2014) e94137.

[53] A. Baronchelli, R. Ferrer-i Cancho, R. Pastor-Satorras, N. Chater, M. H. Christiansen, Networks in cognitive science, Trends in cognitive sciences 17 (7) (2013) 348–360.

[54] C. T. Kello, G. D. Brown, R. Ferrer-i Cancho, J. G. Holden, K. Linkenkaer-Hansen, T. Rhodes, G. C. Van Orden, Scaling laws in cognitive sciences, Trends in cognitive sciences 14 (5) (2010) 223–232.

# 6. Supplementary Information

Table 5: Books used in task of distinguishing three different thematic classes of texts.

| Book | Author | Class |
|---|---|---|
| The Children of the New Forest | Frederick Marryat | C |
| Chronicles of Avonlea | Anne of Green Gables Followed | C |
| Clover | Susan Coolidge | C |
| Eight Cousins | Louisa May Alcott | C |
| The Flying Girl | L. Frank Baum | C |
| A Girl in Ten Thousand | L. T. Meade | C |
| Jack and Jill | Louisa May Alcott | C |
| Jan of the Windmill: A Story of the Plains | Juliana Horatia Gatty Ewing | C |
| The Jungle Book | Rudyard Kipling | C |
| Kim | Rudyard Kipling | C |
| Little Lord Fauntleroy | Frances Hodgson Burnett | C |
| A Little Princess | Frances Hodgson Burnett | C |
| The Magic of Oz | L. Frank Baum | C |
| Mary Louise in the Country | L. Frank Baum | C |
| Peter Pan | J. M. Barrie | C |
| The Princess and Curdie | George MacDonald | C |
| Six to Sixteen: A Story for Girls | Juliana Horatia Gatty Ewing | C |
| Steve and the Steam Engine | Sara Ware Bassett | C |
| A Sweet Little Maid | Amy Ella Blanchard | C |
| Try and Trust; Or, Abner Holden's Bound Boy | Jr. Horatio Alger | C |
| Understood Betsy | Dorothy Canfield Fisher | C |

| | | |
|---|---|---|
| The Water-Babies | Charles Kingsley | C |
| What Katy Did Next | Susan Coolidge | C |
| The Wind in the Willows | Kenneth Grahame | C |
| The Young Explorer; Or, Claiming His Fortune | Jr. Horatio Alger | C |
| The Thirty-Nine Steps | John Buchan | I |
| That Affair at Elizabeth | Burton Egbert Stevenson | I |
| The Mysterious Affair at Styles | Agatha Christie | I |
| An African Millionaire: Episodes in the Life of the Illustrious Colonel Clay | Grant Allen | I |
| The Angel of Terror | Edgar Wallace | I |
| The Crystal Stopper | Maurice Leblanc | I |
| Dead Men's Money | J. S. Fletcher | I |
| Dead Men Tell No Tales | E. W. Hornung | I |
| The Devil Doctor | Sax Rohmer | I |
| The Golden Scorpion | Sax Rohmer | I |
| Greenmantle | John Buchan | I |
| The Hound of the Baskervilles | Arthur Conan Doyle | I |
| Martin Hewitt, Investigator | Arthur Morrison | I |
| The Middle of Things | J. S. Fletcher | I |
| Murder in the Gunroom | H. Beam Piper | I |
| The Mystery of 31 New Inn | R. Austin Freeman | I |
| The Old Man in the Corner | Baroness Emmuska Orczy | I |
| The Passenger from Calais | Arthur Griffiths | I |
| The Red Thumb Mark | R. Austin Freeman | I |
| The Riddle of the Frozen Flame | Mary E. Hanshew and Thomas W. Hanshew | I |
| The Secret Adversary | Agatha Christie | I |

| | | |
|---|---|---|
| The Secret Agent: A Simple Tale | Joseph Conrad | I |
| The Shadow of the Rope | E. W. Hornung | I |
| The Sign of the Four | Arthur Conan Doyle | I |
| A Strange Disappearance | Anna Katharine Green | I |
| Aesthetical Essays of Friedrich Schiller | Friedrich Schiller | P |
| The Analysis of Mind | Bertrand Russell | P |
| Analysis of Mr. Mill's System of Logic | W. Stebbing | P |
| Autobiography | John Stuart Mill | P |
| Beyond Good and Evil | Friedrich Wilhelm Nietzsche | P |
| Democracy and Education: An Introduction to the Philosophy of Education | Dewey | P |
| Discourse on the Method of Rightly Conducting One's Reason and of Seeking Truth | René Descartes | P |
| An Enquiry Concerning Human Understanding | David Hume | P |
| An Enquiry Concerning the Principles of Morals | David Hume | P |
| Essays on some unsettled Questions of Political Economy | John Stuart Mill | P |
| The Ethics of Aristotle | Aristotle | P |
| Jewish History : An Essay in the Philosophy of History | Simon Dubnow | P |
| Laughter: An Essay on the Meaning of the Comic | Henri Bergson | P |

| | | |
|---|---|---|
| On Liberty | John Stuart Mill | P |
| Mind and Motion and Monism | George John Romanes | P |
| The Philosophy of the Moral Feelings | John Abercrombie | P |
| Philosophy and Religion | Hastings Rashdall | P |
| A Pluralistic Universe | William James | P |
| Politics: A Treatise on Governmen | Aristotle | P |
| The Problems of Philosophy | Bertrand Russell | P |
| Proposed Roads to Freedom | Bertrand Russell | P |
| The Psychology of Nations | G. E. Partridge | P |
| The Prince | Niccolò Machiavelli | P |
| Thus Spake Zarathustra: A Book for All and None | Friedrich Wilhelm Nietzsche | P |
| A Treatise Concerning the Principles of Human Knowledge | George Berkeley | P |

Table 6: Books used in task of discriminating books according to their publication dates.

| Book | Author | Year |
|---|---|---|
| The Spectator | Joseph Addison and Sir Richard Steele | 1711 |
| Robinson Crusoe | Daniel Defoe | 1719 |
| A Journal of the Plague Year | Daniel Defoe | 1722 |
| Gulliver's Travels into Several Remote Nations of the World | Jonathan Swift | 1726 |
| Selected Sermons of Jonathan Edwards | Jonathan Edwards | 1731 |
| A Treatise of Human Nature | David Hume | 1738 |

| | | |
|---|---|---|
| Pamela, or Virtue Rewarded | Samuel Richardson | 1740 |
| The Fortunate Foundlings | Eliza Fowler Haywood | 1744 |
| The Adventures of Roderick Random | T. Smollett | 1748 |
| Clarissa Harlowe; or the history of a young lady — Volume 1 | Samuel Richardson | 1748 |
| Life's Progress Through the Passions; Or, The Adventures of Natura | Eliza Fowler Haywood | 1748 |
| History of Tom Jones, a Foundling | Henry Fielding | 1749 |
| Amelia | Henry Fielding | 1751 |
| The Vicar of Wakefield | Oliver Goldsmith | 1761 |
| The Castle of Otranto | Horace Walpole | 1764 |
| The Life and Opinions of Tristram Shandy, Gentleman | Laurence Sterne | 1767 |
| Thoughts on the Present Discontents, and Speeches | Edmund Burke | 1770 |
| The Expedition of Humphry Clinker | T. Smollett | 1771 |
| The Writings of Thomas Paine | Thomas Paine | 1774 |
| A Journey to the Western Islands of Scotland | Samuel Johnson | 1775 |
| Evelina, Or, the History of a Young Lady's Entrance into the World | Fanny Burney | 1778 |
| Cecilia; Or, Memoirs of an Heiress — Volume 1 | Fanny Burney | 1782 |

| | | |
|---|---|---|
| Life of Samuel Johnson | James Boswell | 1791 |
| A Vindication of the Rights of Woman | Mary Wollstonecraft | 1791 |
| A Sicilian Romance | Ann Ward Radcliffe | 1792 |
| The Autobiography of Benjamin Franklin | Benjamin Franklin | 1793 |
| Caleb Williams; Or, Things as They Ar | William Godwin | 1794 |
| The Mysteries of Udolpho | Ann Ward Radcliffe | 1794 |
| Memoirs of Emma Courtney | Mary Hays | 1796 |
| The Monk: A Romanc | M. G. Lewis | 1796 |
| Sense and Sensibility | Jane Austen | 1811 |
| Pride and Prejudice | Jane Austen | 1813 |
| Emma | Jane Austen | 1813 |
| Frankenstein; Or, The Modern Prometheus | Mary Wollstonecraft Shelley | 1818 |
| Persuasion | Jane Austen | 1818 |
| The Spy | James Fenimore Cooper | 1821 |
| The Last of the Mohicans; A narrative of 1757 | James Fenimore Cooper | 1826 |
| The Voyage of the Beagle | Charles Darwin | 1839 |
| The Black Tulip | Alexandre Dumas | 1844 |
| The Three Musketeers | Alexandre Dumas | 1844 |
| Agnes Grey | Anne Brontë | 1847 |
| Wuthering Heights | Emily Brontë | 1847 |
| David Copperfield | Charles Dickens | 1850 |
| The Scarlet Letter | Nathaniel Hawthorne | 1850 |
| The House of the Seven Gables | Nathaniel Hawthorne | 1851 |

| | | |
|---|---|---|
| Moby Dick; Or, The Whale | Herman Melville | 1851 |
| The Boy Hunters | Mayne Reid | 1852 |
| The Confidence-Man: His Masquerade | Herman Melville | 1857 |
| Great Expectations | Charles Dickens | 1861 |
| The Headless Horseman: A Strange Tale of Texas | Mayne Reid | 1865 |
| The Innocents Abroad | Mark Twain | 1869 |
| Daniel Deronda | George Eliot | 1876 |
| Two on a Tower | Thomas Hardy | 1882 |
| Adventures of Huckleberry Finn | Mark Twain | 1884 |
| Stories of Great Men | Faye Huntington | 1887 |
| Life's Little Ironies | Thomas Hardy | 1894 |
| The Time Machine | H. G. Wells | 1895 |
| Dracula | Bram Stoker | 1897 |
| The Invisible Man: A Grotesque Romance | H. G. Wells | 1897 |
| Jane Eyre: An Autobiography | Charlotte Brontë | 1897 |
| Crucial Instances | Edith Wharton | 1901 |
| The Inheritors | Joseph Conrad and Ford Madox Ford | 1901 |
| The Jewel of Seven Stars | Bram Stoker | 1903 |
| The House of Mirth | Edith Wharton | 1905 |
| The Jungle | Upton Sinclair | 1906 |
| The Secret Agent: A Simple Tale | Joseph Conrad | 1907 |
| The Lair of the White Worm | Bram Stoker | 1911 |
| Under Western Eyes | Joseph Conrad | 1911 |
| Daddy-Long-Legs | Jean Webster | 1912 |
| The Lost World | Arthur Conan Doyle | 1912 |

| | | |
|---|---|---|
| The New Freedom | Woodrow Wilson | 1913 |
| Pollyanna | Eleanor H. Porter | 1913 |
| Dubliners | James Joyce | 1914 |
| The Valley of Fear | Arthur Conan Doyle | 1914 |
| Dear Enemy | Jean Webster | 1915 |
| The Good Soldier | Ford Madox Ford | 1915 |
| The Rainbow | D. H. Lawrence | 1915 |
| The Voyage Out | Virginia Woolf | 1915 |
| King Coal : a Novel | Upton Sinclair | 1917 |
| Oh, Money! Money! A Novel | Eleanor H. Porter | 1918 |
| My Antonia | Willa Cather | 1918 |
| Night and Day | Virginia Woolf | 1919 |
| The Age of Innocence | Edith Wharton | 1920 |
| Women in Love | D. H. Lawrence | 1920 |
| This Side of Paradise | F. Scott Fitzgerald | 1920 |
| The Beautiful and Damned | F. Scott Fitzgerald | 1922 |
| The Glimpses of the Moon | Edith Wharton | 1922 |
| One of Ours | Willa Cather | 1922 |
| Ulysses | James Joyce | 1922 |
| The Trial | Franz Kafka | 1925 |