

NLTK através de exemplos:

- Modelagem de tópicos (*topic modeling*)
- Redes de co-ocorrência

Prof. Jesús P. Mena-Chalco
jesus.mena@ufabc.edu.br

2Q-2019

Seeking Life's Bare (Genetic) Necessities

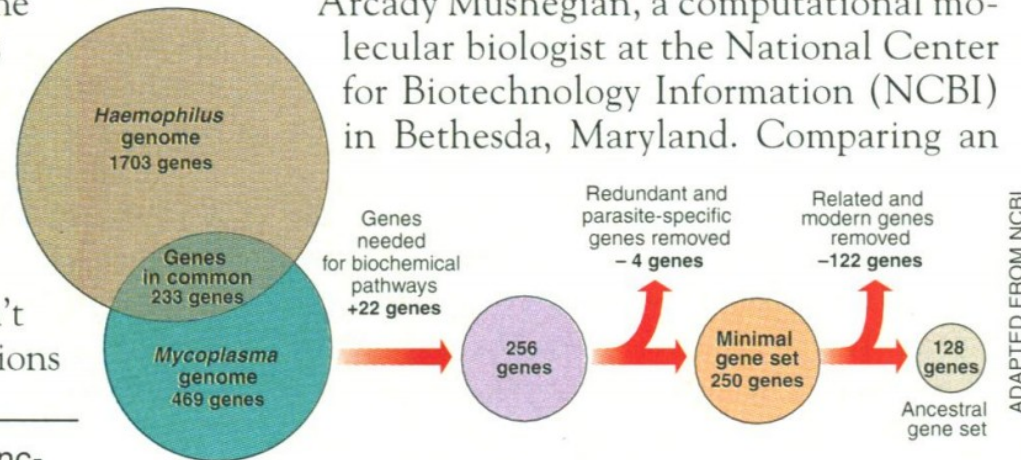
Genética?

COLD SPRING HARBOR, NEW YORK—How many **genes** does an organism need to survive? Last week at the genome meeting here,* two genome researchers with radically different approaches presented complementary views of the basic genes needed for life. One research team, using computer analyses to compare known **genomes**, concluded that today's organisms can be sustained with just 250 genes, and that the earliest life forms required a mere 128 **genes**. The other researcher mapped genes in a simple parasite and estimated that for this organism, 800 genes are plenty to do the job—but that anything short of 100 wouldn't be enough.

Although the numbers don't match precisely, those predictions

* Genome Mapping and Sequencing, Cold Spring Harbor, New York, May 8 to 12.

“are not all that far apart,” especially in comparison to the 75,000 **genes** in the human genome, notes Siv Andersson of Uppsala University in Sweden, who arrived at the 800 number. But coming up with a consensus answer may be more than just a **genetic** numbers game, particularly as more and more **genomes** are completely mapped and sequenced. “It may be a way of organizing any newly **sequenced genome**,” explains Arcady Mushegian, a computational molecular biologist at the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Comparing an



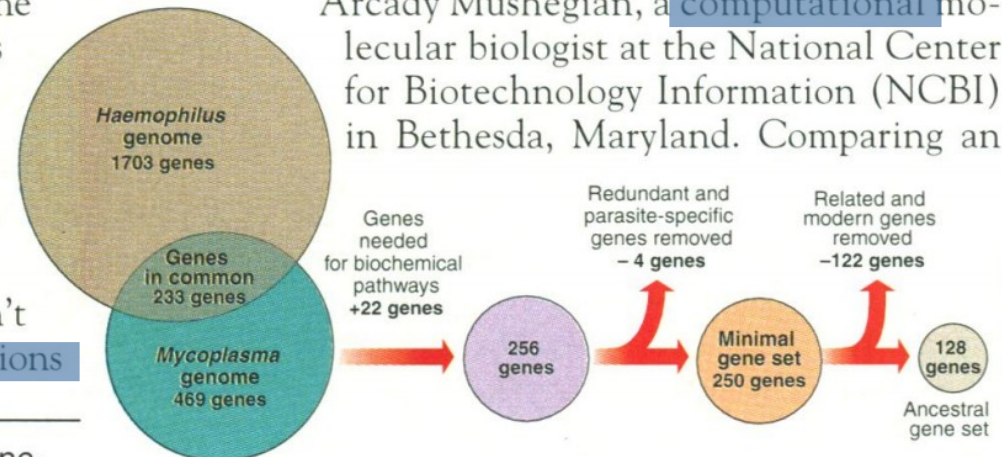
Stripping down. Computer analysis yields an estimate of the minimum modern and ancient genomes.

Seeking Life's Bare (Genetic) Necessities

COLD SPRING HARBOR, NEW YORK—How many **genes** does an organism need to survive? Last week at the genome meeting here,* two genome researchers with radically different approaches presented complementary views of the basic genes needed for life. One research team, using **computer** analyses to compare known **genomes**, concluded that today's organisms can be sustained with just 250 genes, and that the earliest life forms required a mere 128 **genes**. The other researcher mapped genes in a simple parasite and estimated that for this organism, 800 genes are plenty to do the job—but that anything short of 100 wouldn't be enough.

Although the numbers don't match precisely, those **predictions**

“are not all that far apart,” especially in comparison to the 75,000 **genes** in the human genome, notes Siv Andersson of Uppsala University in Sweden, who arrived at the 800 number. But coming up with a consensus answer may be more than just a **genetic numbers** game, particularly as more and more **genomes** are completely mapped and sequenced. “It may be a way of organizing any newly **sequenced genome**,” explains Arcady Mushegian, a **computational** molecular biologist at the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Comparing an



Stripping down. **Computer analysis** yields an estimate of the minimum modern and ancient genomes.

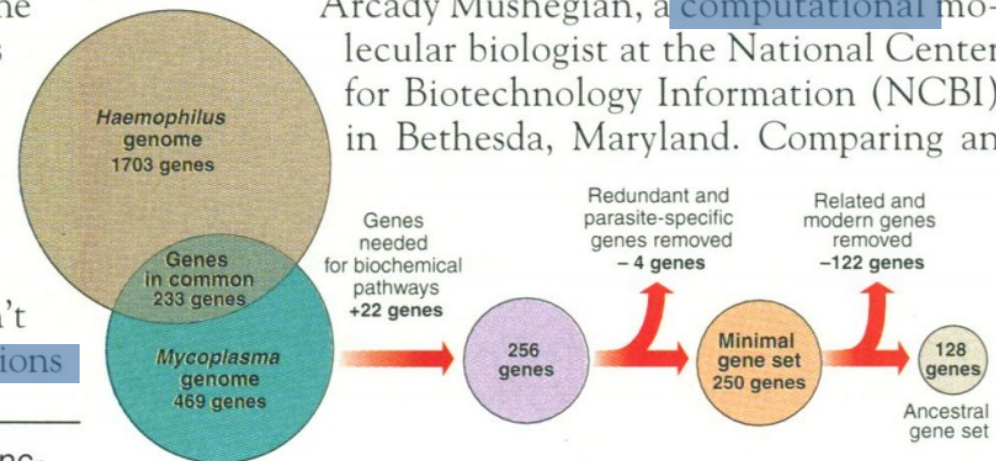
* Genome Mapping and Sequencing, Cold Spring Harbor, New York, May 8 to 12.

Seeking Life's Bare (Genetic) Necessities

COLD SPRING HARBOR, NEW YORK—How many **genes** does an **organism** need to **survive**? Last week at the genome meeting here,* two genome researchers with radically different approaches presented complementary views of the basic genes needed for **life**. One research team, using **computer** analyses to compare known **genomes**, concluded that today's **organisms** can be sustained with just 250 genes, and that the earliest life forms required a mere 128 **genes**. The other researcher mapped genes in a simple parasite and estimated that for this organism, 800 genes are plenty to do the job—but that anything short of 100 wouldn't be enough.

Although the numbers don't match precisely, those **predictions**

“are not all that far apart,” especially in comparison to the 75,000 **genes** in the human genome, notes Siv Andersson of Uppsala University in Sweden, who arrived at the 800 number. But coming up with a consensus answer may be more than just a **genetic numbers** game, particularly as more and more **genomes** are completely mapped and sequenced. “It may be a way of organizing any newly **sequenced genome**,” explains Arcady Mushegian, a **computational** molecular biologist at the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Comparing an



Stripping down. **Computer analysis** yields an estimate of the minimum modern and ancient genomes.

* Genome Mapping and Sequencing, Cold Spring Harbor, New York, May 8 to 12.

Genética?

Computação?

Biologia?

Este exemplo evidencia que um artigo/documento pode estar composto de **diferentes tópicos** (ou unidades) que se misturam para a elaboração de ideias.

Genética?

Computação?

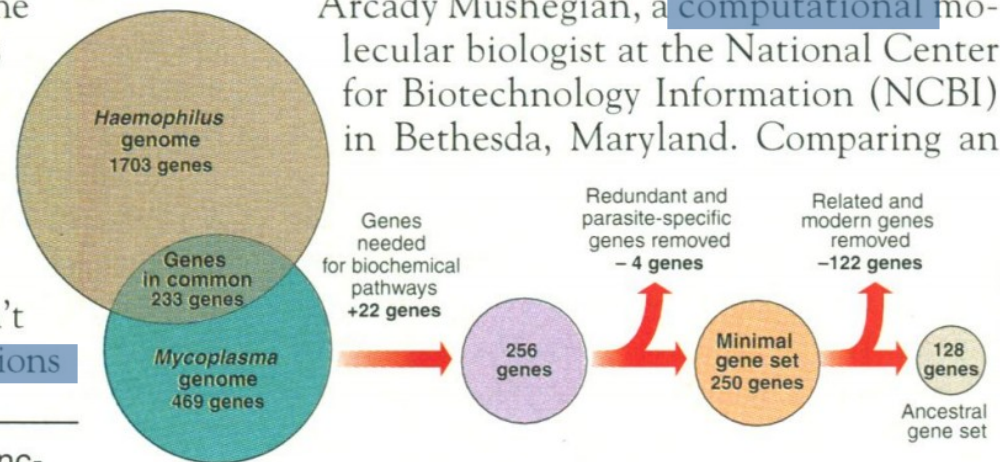
Biologia?

COLD SPRING HARBOR, NEW YORK—How many **genes** does an **organism** need to **survive**? Last week at the genome meeting here,* two genome researchers with radically different approaches presented complementary views of the basic genes needed for **life**. One research team, using **computer** analyses to compare known **genomes**, concluded that today's **organisms** can be sustained with just 250 genes, and that the earliest life forms required a mere 128 **genes**. The other researcher mapped genes in a simple parasite and estimated that for this organism, 800 genes are plenty to do the job—but that anything short of 100 wouldn't be enough.

Although the numbers don't match precisely, those **predictions**

* Genome Mapping and Sequencing, Cold Spring Harbor, New York, May 8 to 12.

“are not all that far apart,” especially in comparison to the 75,000 **genes** in the human genome, notes Siv Andersson of Uppsala University in Sweden, who arrived at the 800 number. But coming up with a consensus answer may be more than just a **genetic numbers** game, particularly as more and more **genomes** are completely mapped and sequenced. “It may be a way of organizing any newly **sequenced genome**,” explains Arcady Mushegian, a **computational** molecular biologist at the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Comparing an



Stripping down. **Computer analysis** yields an estimate of the minimum modern and ancient genomes.

Artigo de revisão em Communication of ACM - 2012

review articles

DOI:10.1145/2133806.2133828

Surveying a suite of algorithms that offer a solution to managing large document archives.

BY DAVID M. BLEI

Probabilistic Topic Models

AS OUR COLLECTIVE knowledge continues to be digitized and stored—in the form of news, blogs, Web pages, scientific articles, books, images, sound, video, and social networks—it becomes more difficult to find and discover what we are looking for. We need new computational tools to help organize, search, and understand these vast amounts of information.

Right now, we work with online information using two main tools—search and links. We type keywords into a search engine and find a set of documents related to them. We look at the documents in that set, possibly navigating to other linked documents. This is a powerful way of interacting with our online archive, but something is missing.

Imagine searching and exploring documents based on the themes that run through them. We might “zoom in” and “zoom out” to find specific or broader themes; we might look at how those themes changed through time or how they are connected to each other. Rather than finding documents through keyword search alone, we might first find the theme that we are interested in, and then examine the documents related to that theme.

For example, consider using themes to explore the complete history of the New York Times. At a broad level, some of the themes might correspond to the sections of the newspaper—foreign policy, national affairs, sports. We could zoom in on a theme of interest, such as foreign policy, to reveal various aspects of it—Chinese foreign policy, the conflict in the Middle East, the U.S.’s relationship with Russia. We could then navigate through time to reveal how these specific themes have changed, tracking, for example, the changes in the conflict in the Middle East over the last 50 years. And, in all of this exploration, we would be pointed to the original articles relevant to the themes. The thematic structure would be a new kind of window through which to explore and digest the collection.

But we do not interact with electronic archives in this way. While more and more texts are available online, we simply do not have the human power to read and study them to provide the kind of browsing experience described above. To this end, machine learning researchers have developed *probabilistic topic modeling*, a suite of algorithms that aim to discover and annotate large archives of documents with thematic information. Topic modeling algorithms are statistical methods that analyze the words of the original texts to discover the themes that run through them, how those themes are connected to each other, and how they change over

» key insights

- Topic models are algorithms for discovering the main themes that pervade a large and otherwise unstructured collection of documents. Topic models can organize the collection according to the discovered themes.
- Topic modeling algorithms can be applied to massive collections of documents. Recent advances in this field allow us to analyze streaming collections, like you might find from a Web API.
- Topic modeling algorithms can be adapted to many kinds of data. Among other applications, they have been used to find patterns in genetic data, images, and social networks.

APRIL 2012 | VOL. 55 | NO. 4 | COMMUNICATIONS OF THE ACM 77

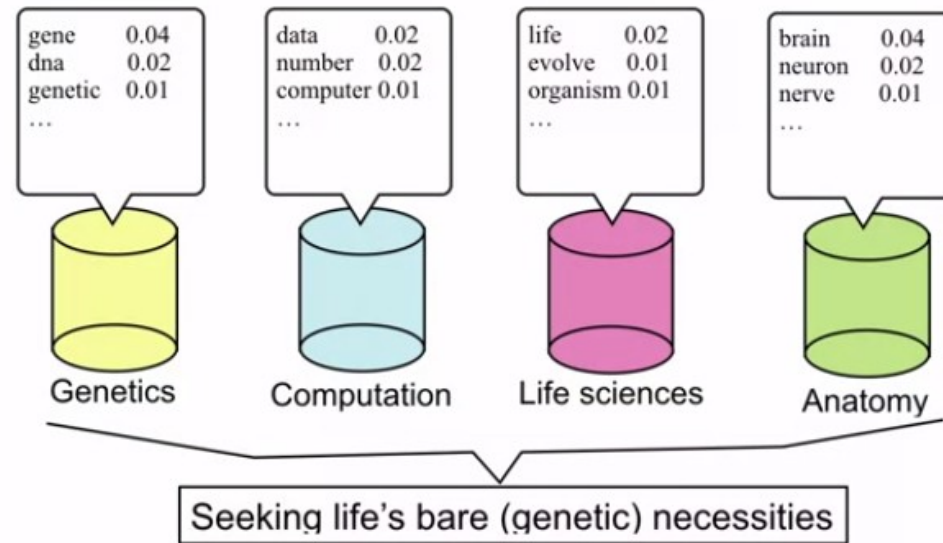
Probabilistic Topic Models | April 2012 | Communications of the ACM

<https://m-cacm.acm.org/magazines/2012/4/...probabilistic-topic-models/fulltext?...true> ▼

by DM Blei - Cited by 4531 - Related articles

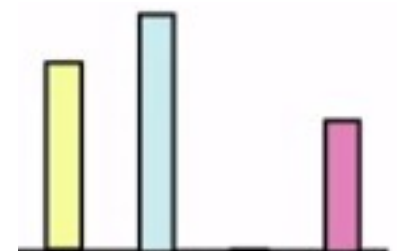
LDA and probabilistic models. LDA and other **topic models** are part of the larger field of **probabilistic modeling**. In generative **probabilistic modeling**, we treat our data as arising from a generative process that includes hidden variables.

Ideia: Um documento textual é uma mistura de tópicos.



Ao termos um ***novo documento*** ao ser analisado, e considerando apenas esses quatro tópicos, podemos pensar que para esse documento, o mais provável é que seja de **computação**.

Novo documento



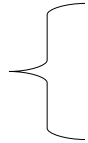
O que é um tópico?

Um tópico pode ser entendido como: **Um tema (ou assunto) de um discurso.**

Os tópicos podem ser representados por uma distribuição de palavras. Isso significa que uma palavra tem certa probabilidade de aparecer nesse tópico.

Exemplo de tópicos

Tópicos



“Genetics”

human
genome
dna
genetic
genes
sequence
gene
molecular
sequencing
map
information
genetics
mapping
project
sequences

“Evolution”

evolution
evolutionary
species
organisms
life
origin
biology
groups
phylogenetic
living
diversity
group
new
two
common

“Disease”

disease
host
bacteria
diseases
resistance
bacterial
new
strains
control
infectious
malaria
parasite
parasites
united
tuberculosis

“Computers”

computer
models
information
data
computers
system
network
systems
model
parallel
methods
networks
software
new
simulations

Modelagem de tópicos (MT)?

É uma análise “abstrata” (em alto nível) do conteúdo de documentos.

É apropriado quando, frente a um **corpus grande**, deseja-se entender, de forma rápida, o conteúdo dele:

- Trata-se de um documento de "esportes"?
- Trata-se de um documento de “genética”?
- Trata-se de um documento de “computação”?

A **MT** também pode ser utilizada para identificar: (i) como esses tópicos se conectam, e (ii) como mudam ao longo do tempo.

Modelagem de tópicos (MT)

O que geralmente é considerado como entrada para a MT são:

- Uma **coleção de documentos** (ou corpus).
- **Número** finito de tópicos.

O que não sabemos:

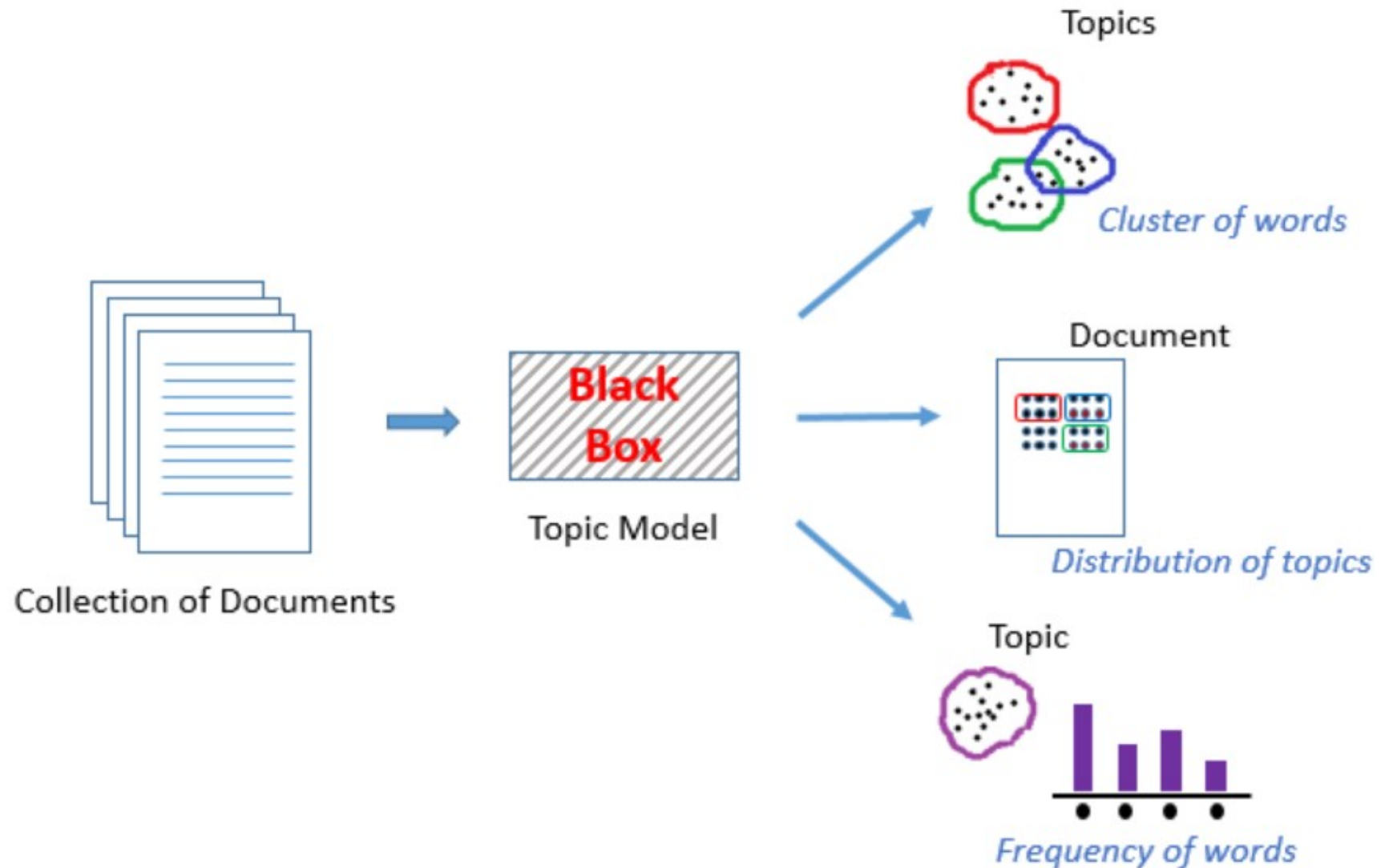
- Os **nomes** dos tópicos.

Não podemos informar, por exemplo, a busca de tópicos em computação.

- A **distribuição de tópicos** para cada documento.

Não sabemos se os documentos tem um conteúdo associado a 50% computação e 50% esportes.

Modelagem de tópicos (MT)



Modelagem de tópicos (MT)

Note, então, que estamos tratando um problema de **agrupamento de textos**, com a variante de que:

Documentos e palavras são agrupadas simultaneamente.

Existem diferentes métodos para modelagem de tópicos, as duas abordagens mais conhecidas são:

- **PLSA:** Probabilistic Latent Semantic Analysis (1999)
- **LDA:** Latent Dirichlet Allocation (2003)

(1) Trabalho pioneiro de Hofmann

289

Probabilistic Latent Semantic Analysis

Thomas Hofmann

EECS Department, Computer Science Division, University of California, Berkeley &
International Computer Science Institute, Berkeley, CA
hofmann@cs.berkeley.edu

Abstract

Probabilistic Latent Semantic Analysis is a novel statistical technique for the analysis of two-mode and co-occurrence data, which has applications in information retrieval and filtering, natural language processing, machine learning from text, and in related areas. Compared to standard Latent Semantic Analysis which stems from linear algebra and performs a Singular Value Decomposition of co-occurrence tables, the proposed method is based on a mixture decomposition derived from a latent class model. This results in a more principled approach which has a solid foundation in statistics. In order to avoid overfitting, we propose a widely applicable generalization of maximum likelihood model fitting by tempered EM. Our approach yields substantial and consistent improvements over Latent Semantic Analysis in a number of experiments.

1 Introduction

Learning from text and natural language is one of the great challenges of Artificial Intelligence and Machine Learning. Any substantial progress in this domain has strong impact on many applications ranging from information retrieval, information filtering, and intelligent interfaces, to speech recognition, natural language processing, and machine translation. One of the fundamental problems is to learn the *meaning* and *usage* of words in a data-driven fashion, i.e., from some given text corpus, possibly without further linguistic prior knowledge.

The main challenge a machine learning system has to address roots in the distinction between the lexical level of “what actually has been said or written” and the semantical level of “what was intended” or “what

was referred to” in a text or an utterance. The resulting problems are twofold: (i) polysems, i.e., a word may have multiple senses and multiple types of usage in different context, and (ii) synonyms and semantically related words, i.e., different words may have a similar meaning, they may at least in certain contexts denote the same concept or – in a weaker sense – refer to the same topic.

Latent semantic analysis (LSA) [3] is well-known technique which partially addresses these questions. The key idea is to map high-dimensional count vectors, such as the ones arising in vector space representations of text documents [12], to a lower dimensional representation in a so-called *latent semantic space*. As the name suggests, the goal of LSA is to find a data mapping which provides information well beyond the lexical level and reveals semantical relations between the entities of interest. Due to its generality, LSA has proven to be a valuable analysis tool with a wide range of applications (e.g. [3, 5, 8, 1]). Yet its theoretical foundation remains to a large extent unsatisfactory and incomplete.

This paper presents a statistical view on LSA which leads to a new model called *Probabilistic Latent Semantics Analysis* (PLSA). In contrast to standard LSA, its probabilistic variant has a sound statistical foundation and defines a proper generative model of the data. A detailed discussion of the numerous advantages of PLSA can be found in subsequent sections.

2 Latent Semantic Analysis

2.1 Count Data and Co-occurrence Tables

LSA can in principle be applied to any type of count data over a discrete dyadic domain (cf. [7]). However, since the most prominent application of LSA is in the analysis and retrieval of text documents, we focus on this setting for sake of concreteness. Suppose therefore we have given a collection of text doc-

HOFMANN, Thomas.

Probabilistic latent semantic analysis.

In: Proceedings of the Fifteenth conference on Uncertainty in artificial intelligence.

Morgan Kaufmann Publishers Inc., 1999. p. 289-296.

Citado por 2538

(2) Trabalho mais popular de Blei *et al.*

Journal of Machine Learning Research 3 (2003) 993-1022

Submitted 2/02; Published 1/03

Latent Dirichlet Allocation

David M. Blei

Computer Science Division
University of California
Berkeley, CA 94720, USA

BLEI@CS.BERKELEY.EDU

Andrew Y. Ng

Computer Science Department
Stanford University
Stanford, CA 94305, USA

ANG@CS.STANFORD.EDU

Michael I. Jordan

Computer Science Division and Department of Statistics
University of California
Berkeley, CA 94720, USA

JORDAN@CS.BERKELEY.EDU

Editor: John Lafferty

Abstract

We describe *latent Dirichlet allocation* (LDA), a generative probabilistic model for collections of discrete data such as text corpora. LDA is a three-level hierarchical Bayesian model, in which each item of a collection is modeled as a finite mixture over an underlying set of topics. Each topic is, in turn, modeled as an infinite mixture over an underlying set of topic probabilities. In the context of text modeling, the topic probabilities provide an explicit representation of a document. We present efficient approximate inference techniques based on variational methods and an EM algorithm for empirical Bayes parameter estimation. We report results in document modeling, text classification, and collaborative filtering, comparing to a mixture of unigrams model and the probabilistic LSI model.

1. Introduction

In this paper we consider the problem of modeling text corpora and other collections of discrete data. The goal is to find short descriptions of the members of a collection that enable efficient processing of large collections while preserving the essential statistical relationships that are useful for basic tasks such as classification, novelty detection, summarization, and similarity and relevance judgments.

Significant progress has been made on this problem by researchers in the field of information retrieval (IR) (Baeza-Yates and Ribeiro-Neto, 1999). The basic methodology proposed by IR researchers for text corpora—a methodology successfully deployed in modern Internet search engines—reduces each document in the corpus to a vector of real numbers, each of which represents ratios of counts. In the popular *tf-idf* scheme (Salton and McGill, 1983), a basic vocabulary of “words” or “terms” is chosen, and, for each document in the corpus, a count is formed of the number of occurrences of each word. After suitable normalization, this term frequency count is compared to an inverse document frequency count, which measures the number of occurrences of a

©2003 David M. Blei, Andrew Y. Ng and Michael I. Jordan.

BLEI, David M.; NG, Andrew Y.;
JORDAN, Michael I.

Latent dirichlet allocation.

Journal of machine Learning
research

v. 3, n. Jan, p. 993-1022, 2003.

Citado por 27564

“Arts”	“Budgets”	“Children”	“Education”
NEW	MILLION	CHILDREN	SCHOOL
FILM	TAX	WOMEN	STUDENTS
SHOW	PROGRAM	PEOPLE	SCHOOLS
MUSIC	BUDGET	CHILD	EDUCATION
MOVIE	BILLION	YEARS	TEACHERS
PLAY	FEDERAL	FAMILIES	HIGH
MUSICAL	YEAR	WORK	PUBLIC
BEST	SPENDING	PARENTS	TEACHER
ACTOR	NEW	SAYS	BENNETT
FIRST	STATE	FAMILY	MANIGAT
YORK	PLAN	WELFARE	NAMPHY
OPERA	MONEY	MEN	STATE
THEATER	PROGRAMS	PERCENT	PRESIDENT
ACTRESS	GOVERNMENT	CARE	ELEMENTARY
LOVE	CONGRESS	LIFE	HAITI

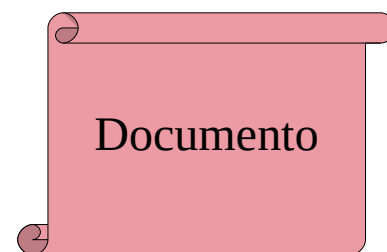
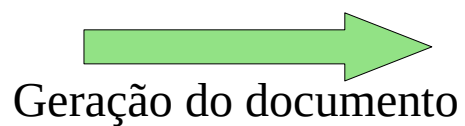
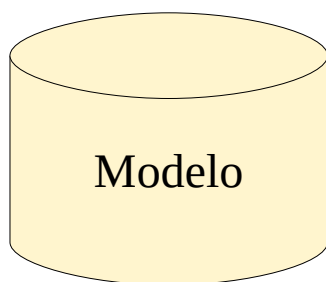
The William Randolph Hearst Foundation will give \$1.25 million to Lincoln Center, Metropolitan Opera Co., New York Philharmonic and Juilliard School. “Our board felt that we had a real opportunity to make a mark on the future of the performing arts with these grants an act every bit as important as our traditional areas of support in health, medical research, education and the social services,” Hearst Foundation President Randolph A. Hearst said Monday in announcing the grants. Lincoln Center’s share will be \$200,000 for its new building, which will house young artists and provide new public facilities. The Metropolitan Opera Co. and New York Philharmonic will receive \$400,000 each. The Juilliard School, where music and the performing arts are taught, will get \$250,000. The Hearst Foundation, a leading supporter of the Lincoln Center Consolidated Corporate Fund, will make its usual annual \$100,000 donation, too.

Figure 8: An example article from the AP corpus. Each color codes a different factor from which the word is putatively generated.

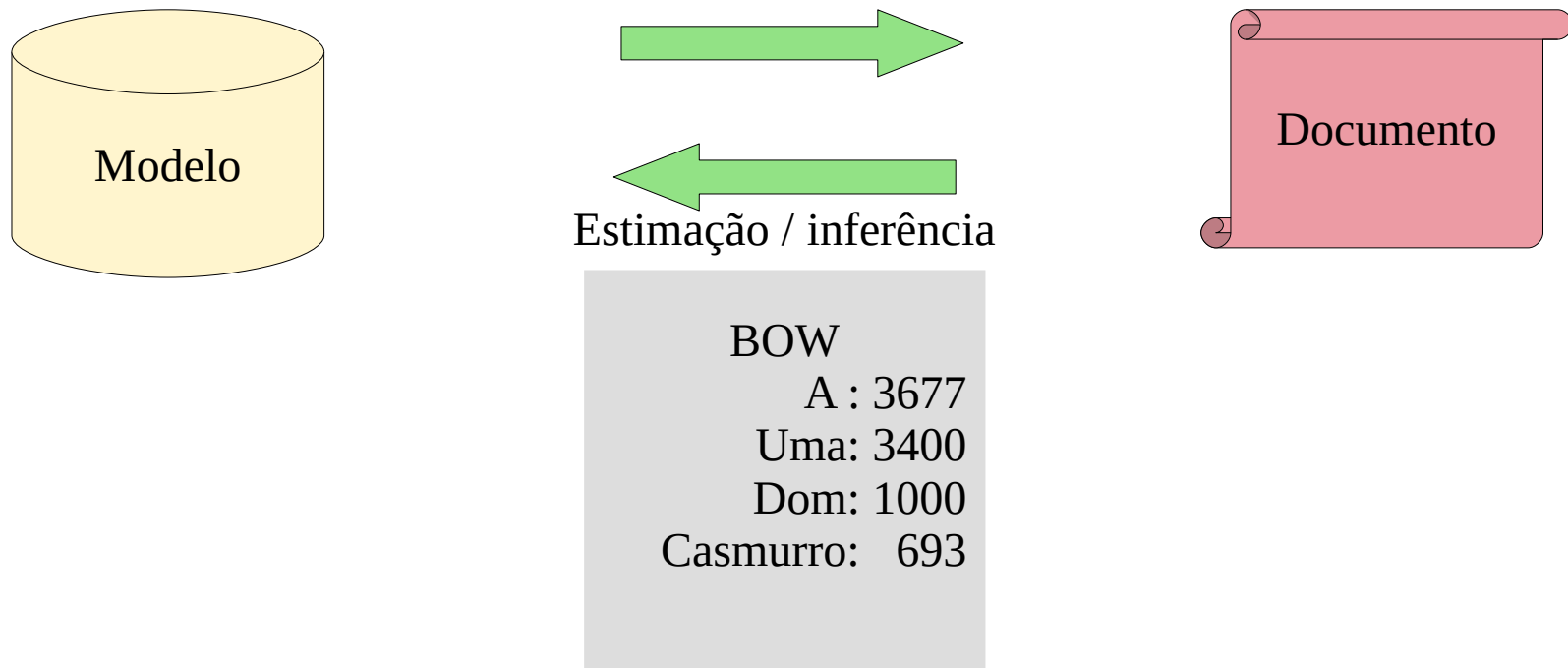


Modelos generativos e LDA

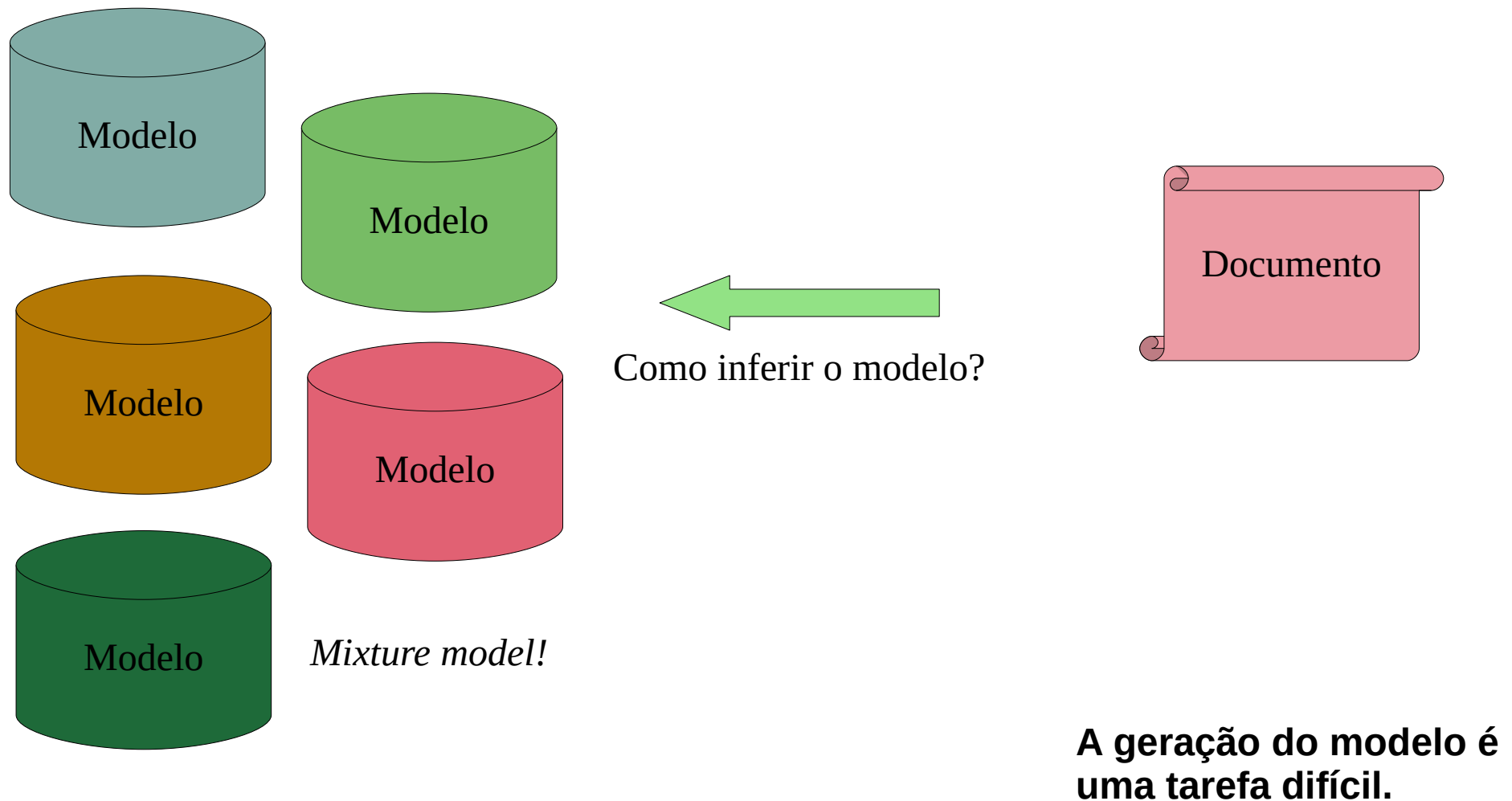
Modelos generativos para textos



Modelos generativos para textos



Modelos generativos para textos



LDA (Latent Dirichlet Allocation)

Topics

gene	0.04
dna	0.02
genetic	0.01
...	

life	0.02
evolve	0.01
organism	0.01
...	

brain	0.04
neuron	0.02
nerve	0.01
...	

data	0.02
number	0.02
computer	0.01
...	

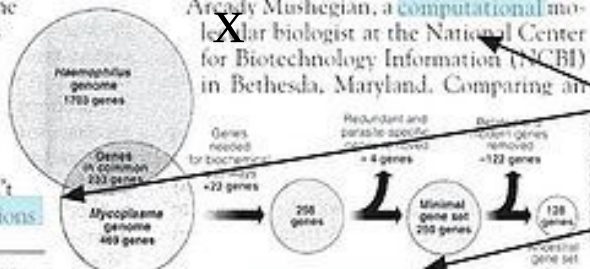
Documents

Seeking Life's Bare (Genetic) Necessities

COLD SPRING HARBOR, NEW YORK—How many **genes** does an **organism** need to **survive**? Last week at the genome meeting here,* two genome researchers with radically different approaches presented complementary views of the basic genes needed for **life**. One research team, using **computer** analyses to compare known **genomes**, concluded that today's **organisms** can be sustained with just 250 genes, and that the earliest life forms required a mere 128 **genes**. The other researcher mapped genes in a simple parasite and estimated that for this organism, 800 genes are plenty to do the job—but that anything short of 100 wouldn't be enough.

Although the numbers don't match precisely, those **predictions**

"are not all that far apart," especially in comparison to the 75,000 **genes** in the human genome, notes Siv Andersson, Uppsala University in Sweden, who arrived at the 800 number. But coming up with a consensus answer may be more than just a **genetic** numbers game, particularly as more and more **genomes** are completely mapped and sequenced. "It may be a way of organizing any newly **sequenced genome**," explains Arcady Mushegian, a **computational** molecular biologist at the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Comparing an

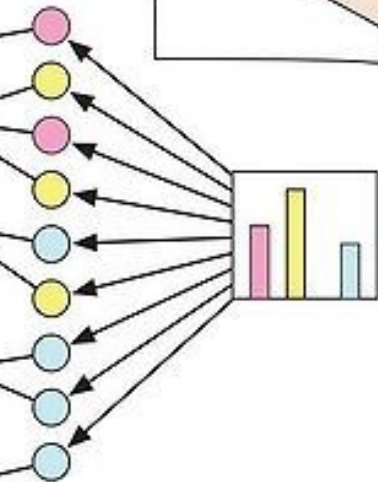


* Genome Mapping and Sequencing, Cold Spring Harbor, New York, May 8 to 12.

Stripping down. Computer analysis yields an estimate of the minimum modern and ancient genomes.

SCIENCE • VOL. 272 • 24 MAY 1996

Topic proportions and assignments



Documento

Tópico A
Tópico B
Tópico C

Tópico A

Palavra 1 (20%)
Palavra 2 (13%)
Palavra 3 (10%)

Tópico C

Palavra 1 (16%)
Palavra 2 (02%)
Palavra 3 (04%)

Tópico B

Palavra 1 (25%)
Palavra 2 (11%)
Palavra 3 (08%)



Atribuir aleatoriamente os tópicos às palavras

A	B	C
B	D	E
A	C	C
C	B	A

(1)

C	E	D
D	A	D
A	B	E
E	C	A

(2)

A	A	B
C	A	B
B	D	E
E	A	B

(3)

E	E	B
B	A	E
D	E	A
C	E	E

(4)

E	A	E
B	E	B
D	A	E
A	D	B

(5)



Calcular $P(T \mid D)$ e $P(W \mid T)$

A	B	C
B	D	E
A	C	C
C	B	A

(1)

C	E	D
D	A	D
A	B	E
E	C	A

(2)

A	A	B
C	A	B
B	D	E
E	A	B

(3)

E	E	B
B	A	E
D	E	A
C	E	E

(4)

E	A	E
B	E	B
D	A	E
A	D	B

(5)

- A (0,05); B(0,00); C(0,05); D(0,07); E(0,10)
- A (0,10); B(0,07); C(0,00); D(0,03); E(0,03)
- A (0,10); B(0,15); C(0,08); D(0,03); E(0,13)

● 0,333
● 0,167
● 0,500

A	B	C
B	D	E
A	C	C
C	B	A

(1)

● 0,250
● 0,250
● 0,500

C	E	D
D	A	D
A	B	E
E	C	A

(2)

● 0,167
● 0,250
● 0,583

A	A	B
C	A	B
B	D	E
E	A	B

(3)

● 0,250
● 0,250
● 0,667

E	E	B
B	A	E
D	E	A
C	E	E

(4)

● 0,333
● 0,250
● 0,417

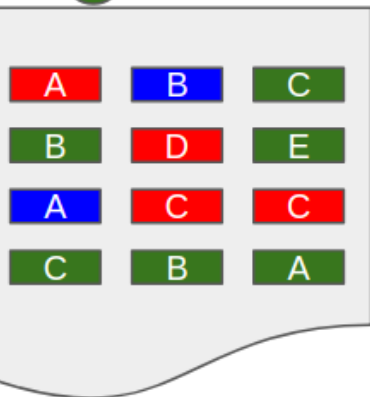
E	A	E
B	E	B
D	A	E
A	D	B

(5)

Atualizar o tópico (T) de cada palavra (W) com $P(T | D) \times P(W | T)$

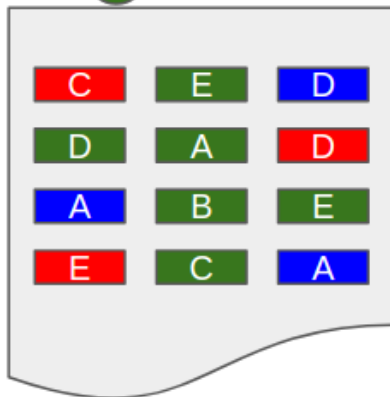
- A (0,05); B(0,00); C(0,05); D(0,07); E(0,10)
- A (0,10); B(0,07); C(0,00); D(0,03); E(0,03)
- A (0,10); B(0,15); C(0,08); D(0,03); E(0,13)

● 0,333
● 0,167
● 0,500



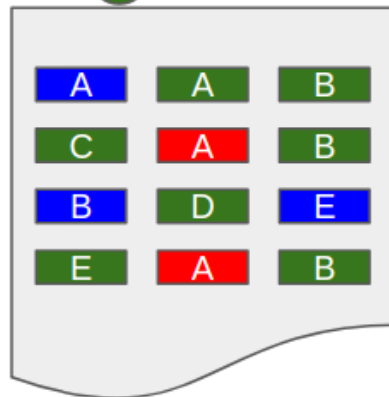
(1)

● 0,250
● 0,250
● 0,500



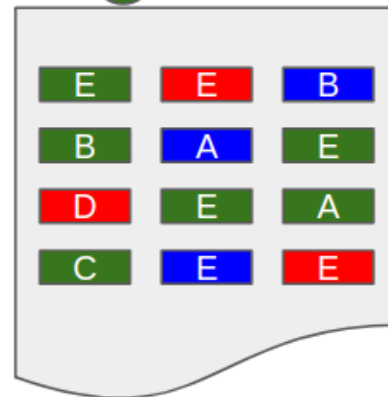
(2)

● 0,167
● 0,250
● 0,583



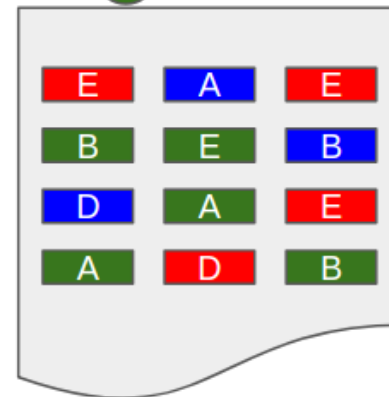
(3)

● 0,250
● 0,250
● 0,667

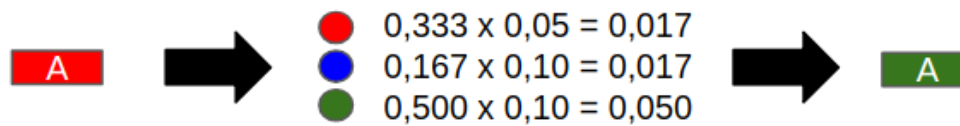


(4)

● 0,333
● 0,250
● 0,417



(5)



- A (0,03); B(0,00); C(0,05); D(0,07); E(0,10)
- A (0,10); B(0,07); C(0,00); D(0,03); E(0,03)
- A (0,12); B(0,15); C(0,08); D(0,03); E(0,13)

● 0,250
● 0,167
● 0,583

A	B	C
B	D	E
A	C	C
C	B	A

(1)

● 0,250
● 0,250
● 0,500

C	E	D
D	A	D
A	B	E
E	C	A

(2)

● 0,167
● 0,250
● 0,583

A	A	B
C	A	B
B	D	E
E	A	B

(3)

● 0,250
● 0,250
● 0,667

E	E	B
B	A	E
D	E	A
C	E	E

(4)

● 0,333
● 0,250
● 0,417

E	A	E
B	E	B
D	A	E
A	D	B

(5)

- A (0,03); B(0,00); C(0,05); D(0,07); E(0,10)
- A (0,10); B(0,07); C(0,00); D(0,03); E(0,03)
- A (0,12); B(0,15); C(0,08); D(0,03); E(0,13)

● 0,250
● 0,167
● 0,583

A	B	C
B	D	E
A	C	C
C	B	A

(1)

● 0,250
● 0,250
● 0,500

C	E	D
D	A	D
A	B	E
E	C	A

(2)

● 0,167
● 0,250
● 0,583

A	A	B
C	A	B
B	D	E
E	A	B

(3)

● 0,250
● 0,250
● 0,667

E	E	B
B	A	E
D	E	A
C	E	E

(4)

● 0,333
● 0,250
● 0,417

E	A	E
B	E	B
D	A	E
A	D	B

(5)

B



● $0,250 \times 0,00 = 0,000$
● $0,167 \times 0,07 = 0,012$
● $0,583 \times 0,15 = 0,087$



B

- A (0,03); B(0,00); C(0,05); D(0,07); E(0,10)
- A (0,10); B(0,05); C(0,00); D(0,03); E(0,03)
- A (0,12); B(0,17); C(0,08); D(0,03); E(0,13)

● 0,250
● 0,083
● 0,667

A	B	C
B	D	E
A	C	C
C	B	A

(1)

● 0,250
● 0,250
● 0,500

C	E	D
D	A	D
A	B	E
E	C	A

(2)

● 0,167
● 0,250
● 0,583

A	A	B
C	A	B
B	D	E
E	A	B

(3)

● 0,250
● 0,250
● 0,667

E	E	B
B	A	E
D	E	A
C	E	E

(4)

● 0,333
● 0,250
● 0,417

E	A	E
B	E	B
D	A	E
A	D	B

(5)

Exemplo (simples)

Suponha termos 5 documentos (com poucas palavras)

- 1) I like to eat broccoli and bananas.
- 2) I ate a banana and spinach smoothie for breakfast.
- 3) Chinchillas and kittens are cute.
- 4) My sister adopted a kitten yesterday.
- 5) **Look at this cute hamster munching on a piece of broccoli.**

Exemplo (texto padronizado)

Sem stop words, e após aplicação de um stemmer:

- 1) like, broccoli, banana
- 2) banana, spinach, smoothi, breakfast
- 3) chinchilla, kitten, cute
- 4) sister, adopt, kitten, yesterday
- 5) look, cute, hamster, munch, piec, broccoli

Exemplo (texto padronizado)

- 1) Topic: $0.117 \cdot \text{"broccoli"} + 0.116 \cdot \text{"banana"} + 0.076 \cdot \text{"cute"} + 0.072 \cdot \text{"look"} + 0.071 \cdot \text{"hamster"}$
- 2) Topic: $0.117 \cdot \text{"broccoli"} + 0.116 \cdot \text{"banana"} + 0.076 \cdot \text{"cute"} + 0.072 \cdot \text{"look"} + 0.071 \cdot \text{"hamster"}$
- 3) Topic: $0.158 \cdot \text{"kitten"} + 0.095 \cdot \text{"sister"} + 0.095 \cdot \text{"yesterday"} + 0.095 \cdot \text{"adopt"} + 0.094 \cdot \text{"chinchilla"}$
- 4) Topic: $0.158 \cdot \text{"kitten"} + 0.095 \cdot \text{"sister"} + 0.095 \cdot \text{"yesterday"} + 0.095 \cdot \text{"adopt"} + 0.094 \cdot \text{"chinchilla"}$
- 5) Topic: $0.117 \cdot \text{"broccoli"} + 0.116 \cdot \text{"banana"} + 0.076 \cdot \text{"cute"} + 0.072 \cdot \text{"look"} + 0.071 \cdot \text{"hamster"}$



Redes de co-ocorrência aplicadas para a caracterização de obras literárias

Universidade Federal do ABC

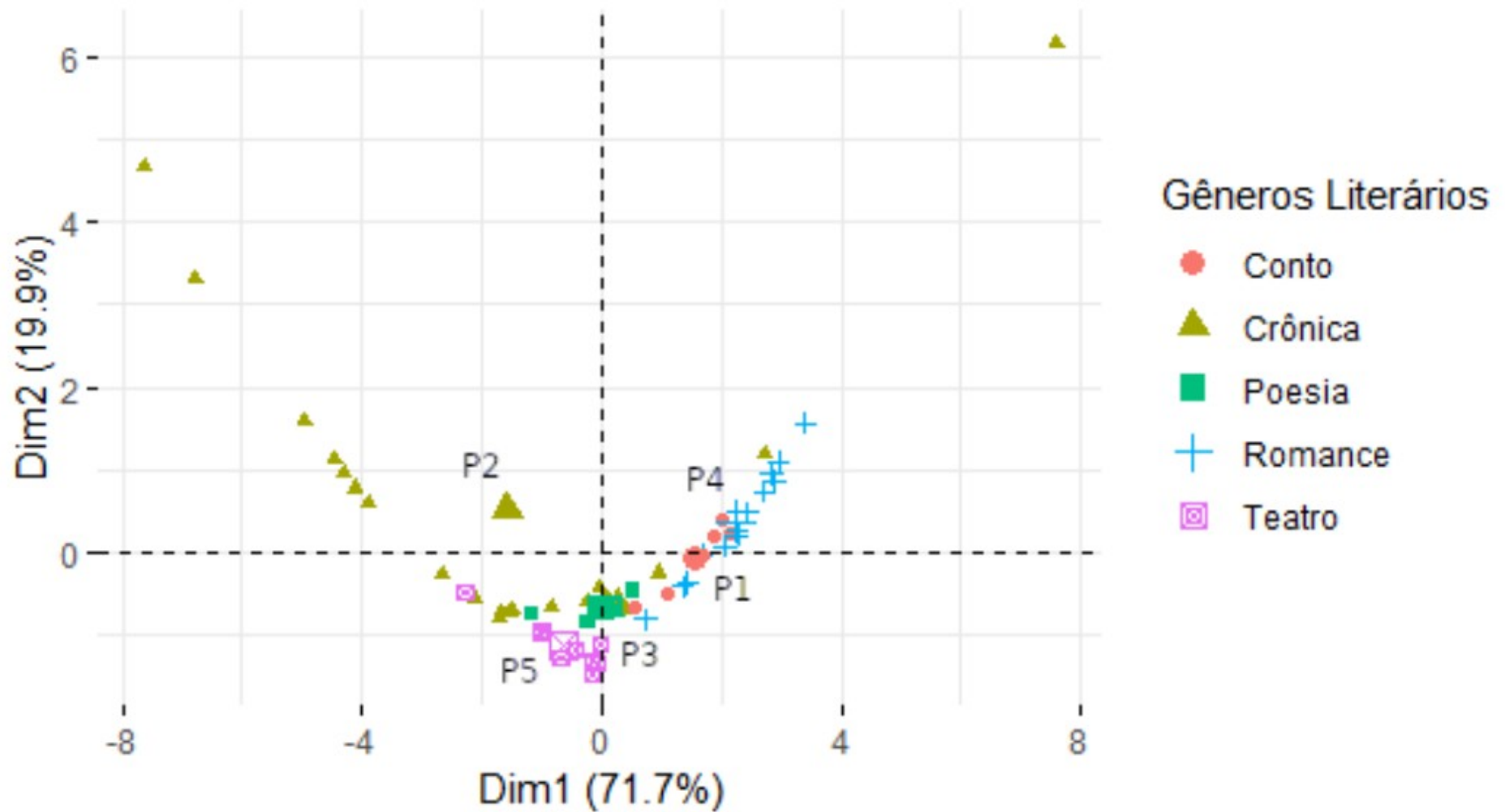
Caracterização de obras literárias usando redes de co-ocorrência

Relatório Final - Iniciação Científica - FAPESP

Processo 2017/26981-9

Aluna: Bruna Pereira Santos

Podemos classificar obras literárias?



Objetivo do trabalho

Demonstrar, através de diferentes testes empíricos, que a **classificação de textos de obras literárias** pode ser realizada usando suas características estruturais (padrões topológicos) extraídas de redes de co-ocorrência textual.

Um grafo/rede de co-ocorrência

❖ **Frase original (Machado de Assis):**

“Talvez porque nenhuma **tinha os** olhos **de** ressaca, **nem os de** cigana oblíqua **e** dissimulada.”

Um grafo/rede de co-ocorrência

- ❖ **Frase original (Machado de Assis):**

“Talvez porque nenhuma **tinha os** olhos **de** ressaca, **nem os de** cigana oblíqua **e** dissimulada.”

- ❖ **Normalização (sem *stemming*):**

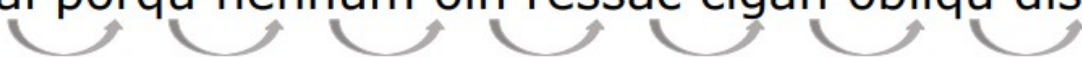
talvez porque nenhuma olhos ressaca cigana oblíqua
dissimulada

- ❖ ***Stemming*:**

tal porqu nenhum olh ressac cigan obliqu dissimul

- ❖ **Rede de co-ocorrência (janela de conexão 2):**

tal porqu nenhum olh ressac cigan obliqu dissimul

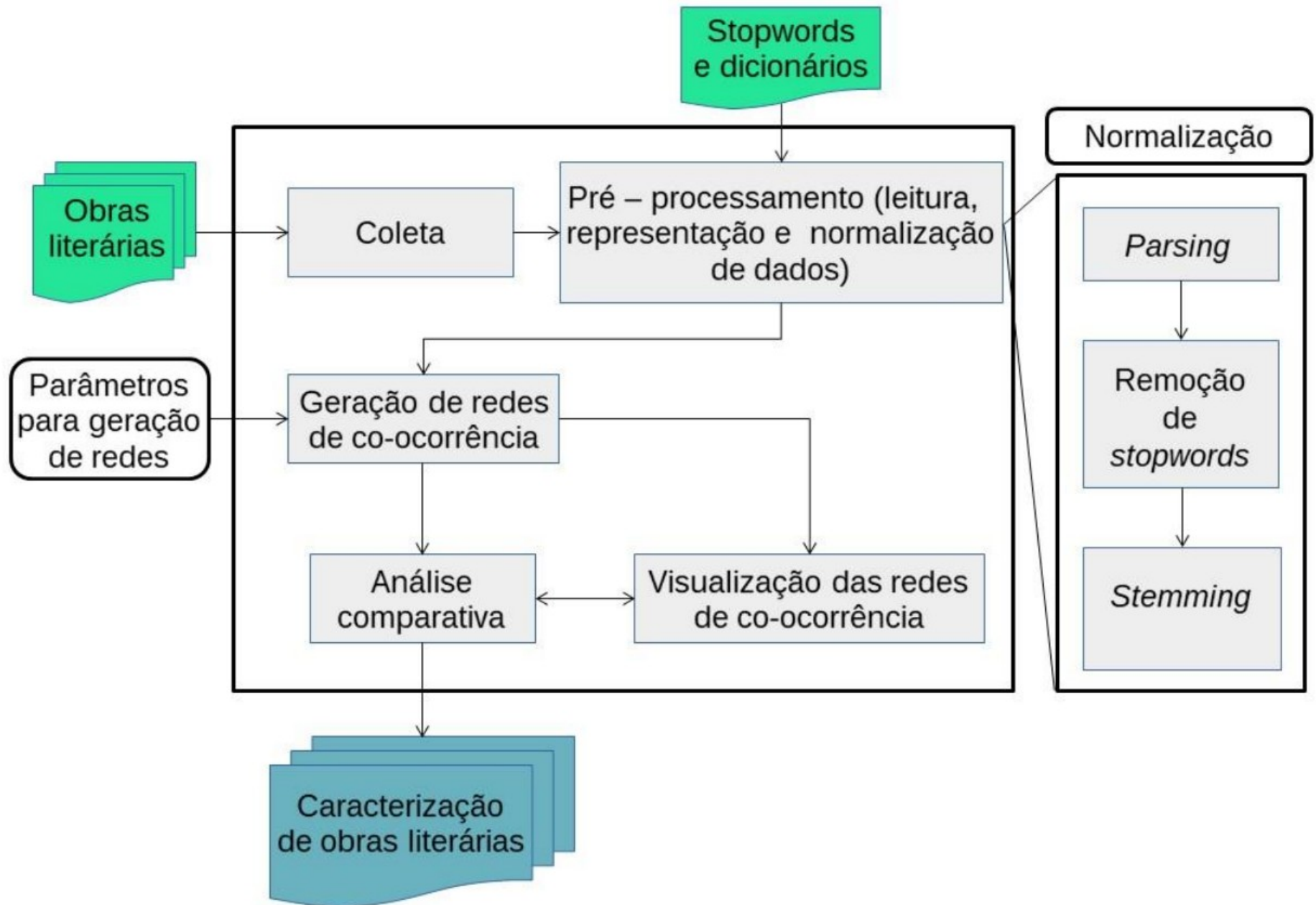


- ❖ **Rede de co-occorrência (janela de conexão 5):**

tal porqu nenhum olh ressac cigan obliqu dissimul



Pipeline



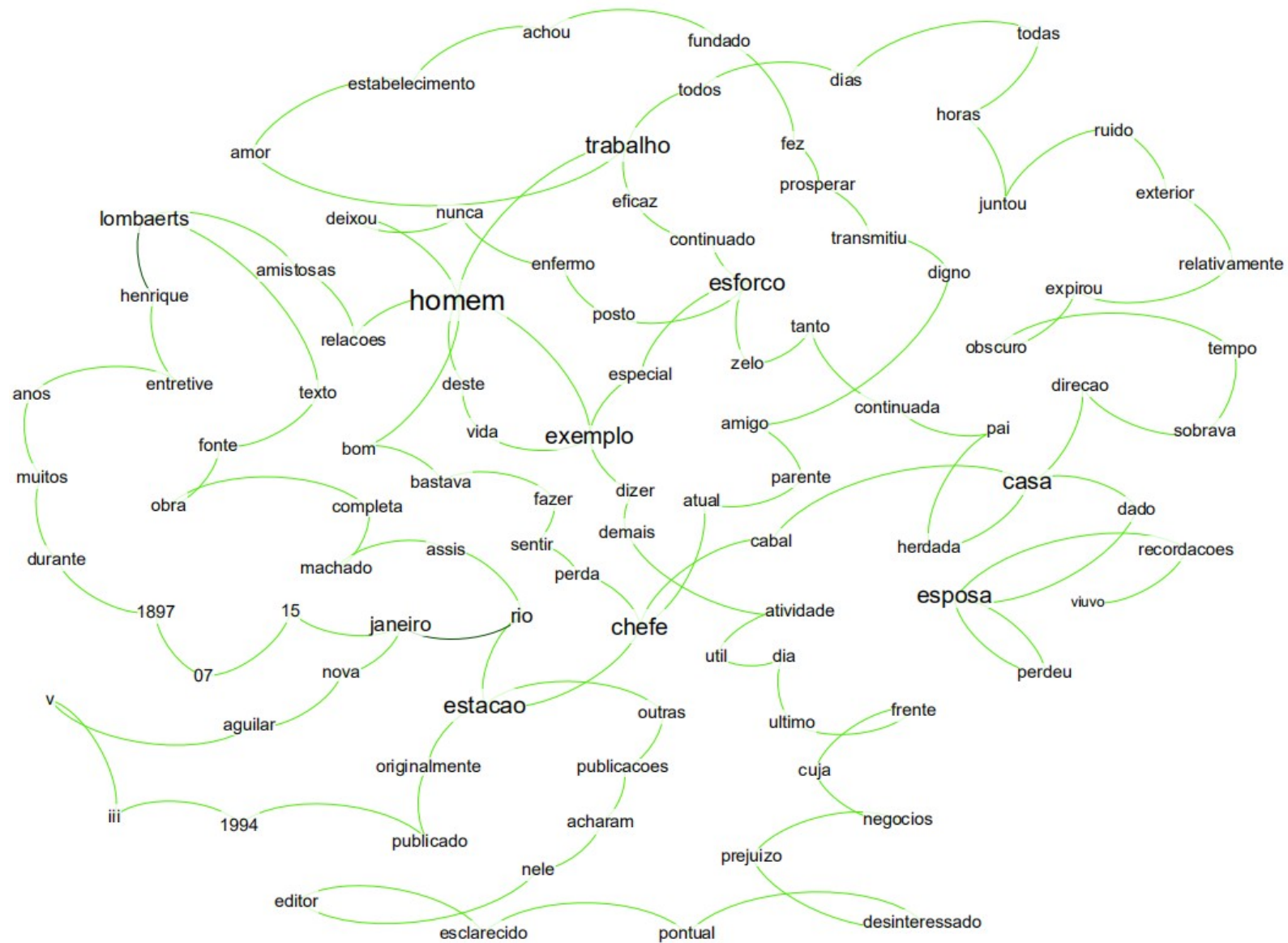
Características: Janela de tamanho 2

Obras/Atributos	Tipo de Grafo	Vértices (vocabulário)	Arestas	Maior frequência entre palavras
Henrique Lobaerts sem <i>stemming</i>	Dirigido	100	110	2
Henrique Lobaerts com <i>stemming</i>	Dirigido	95	110	2
A semana sem <i>stemming</i>	Dirigido	23315	130546	94
A semana com <i>stemming</i>	Dirigido	10484	121179	96
A semana sem <i>stemming</i> (filtro aplicado)	Dirigido	227	185	94
A semana com <i>stemming</i> (filtro aplicado)	Dirigido	237	246	96

Características: Janela de tamanho 5

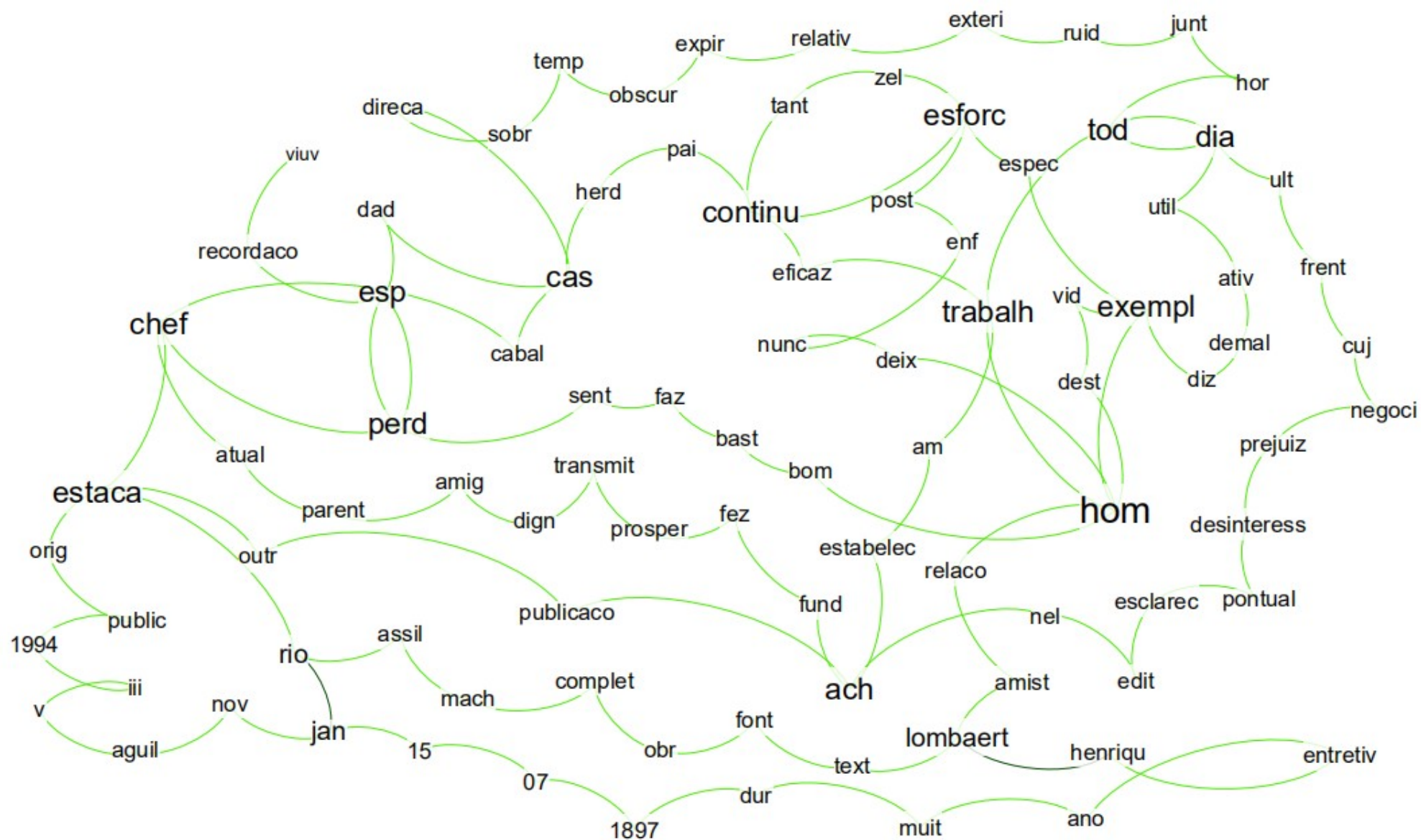
Obras/Atributos	Tipo de Grafo	Vértices (vocabulário)	Arestas	Maior peso das arestas
Henrique Lobaerts sem <i>stemming</i>	Dirigido	100	436	2
Henrique Lobaerts com <i>stemming</i>	Dirigido	95	432	2
A semana sem <i>stemming</i>	Dirigido	23315	517191	100
A semana com <i>stemming</i>	Dirigido	10484	445152	109
A semana sem <i>stemming</i> (filtro aplicado)	Dirigido	279	358	100
A semana com <i>stemming</i> (filtro aplicado)	Dirigido	339	1018	109

Redes de co-ocorrência (J=2)

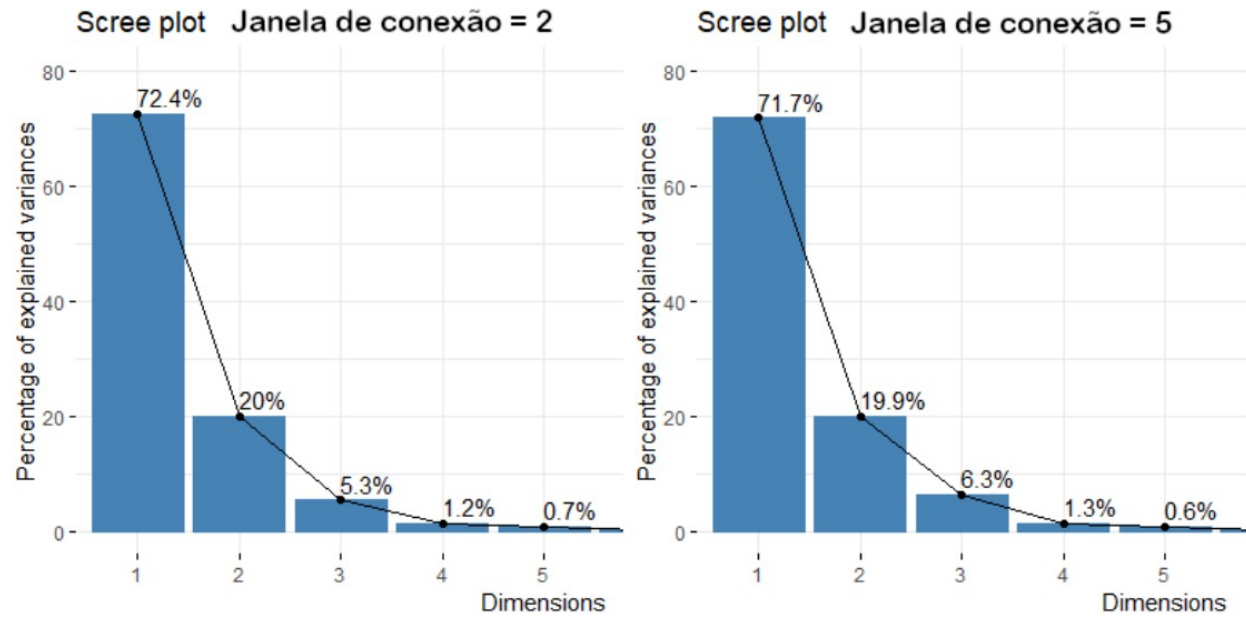


(a) Sem aplicação do *stemming*

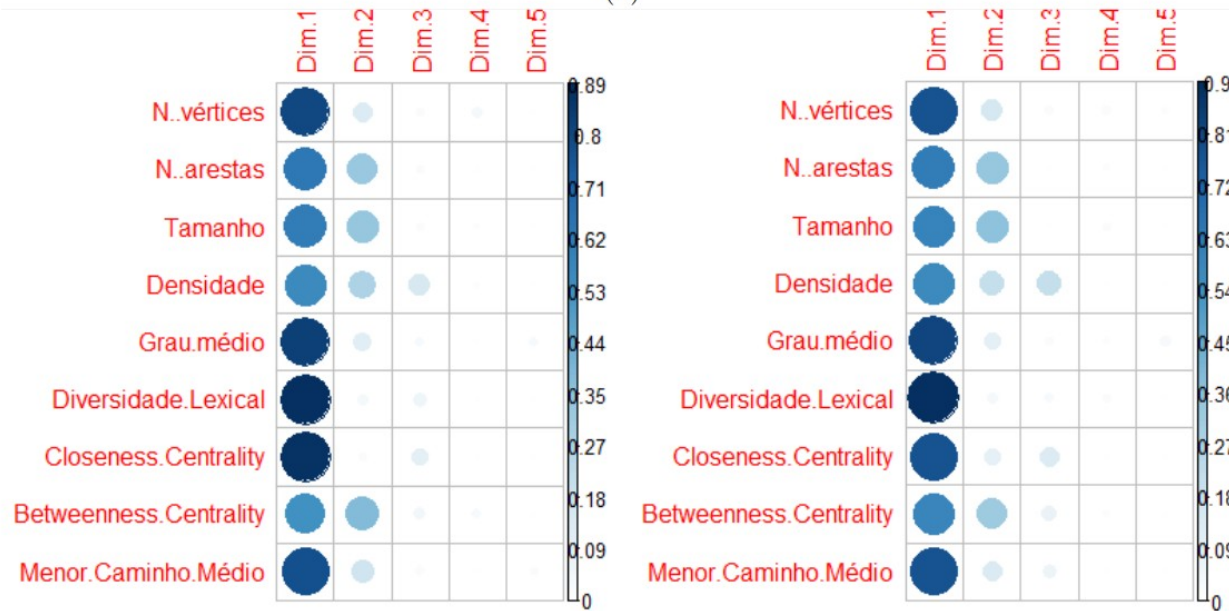
Redes de co-ocorrência (J=2)



(b) Com aplicação do *stemming*

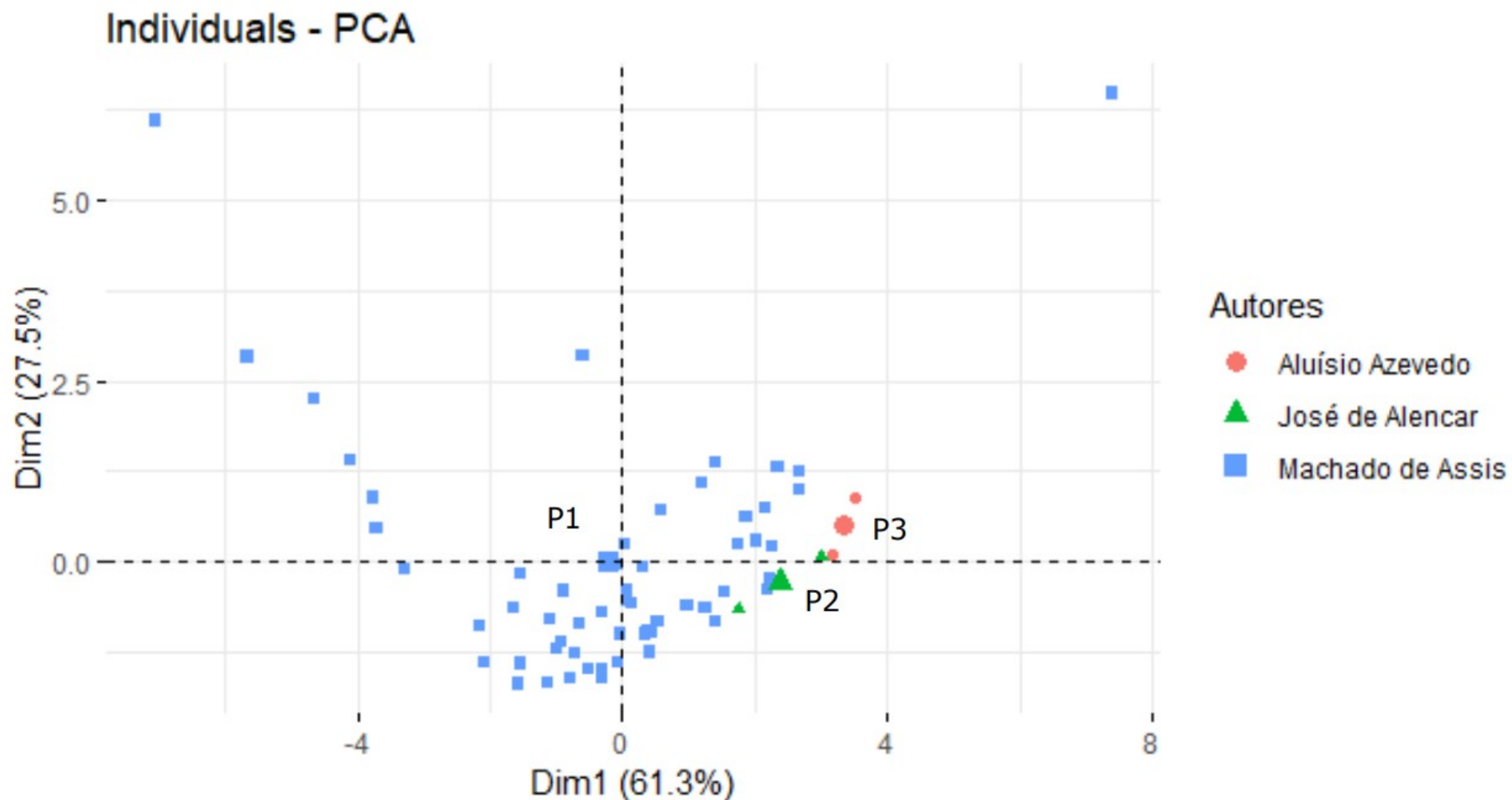


(a)

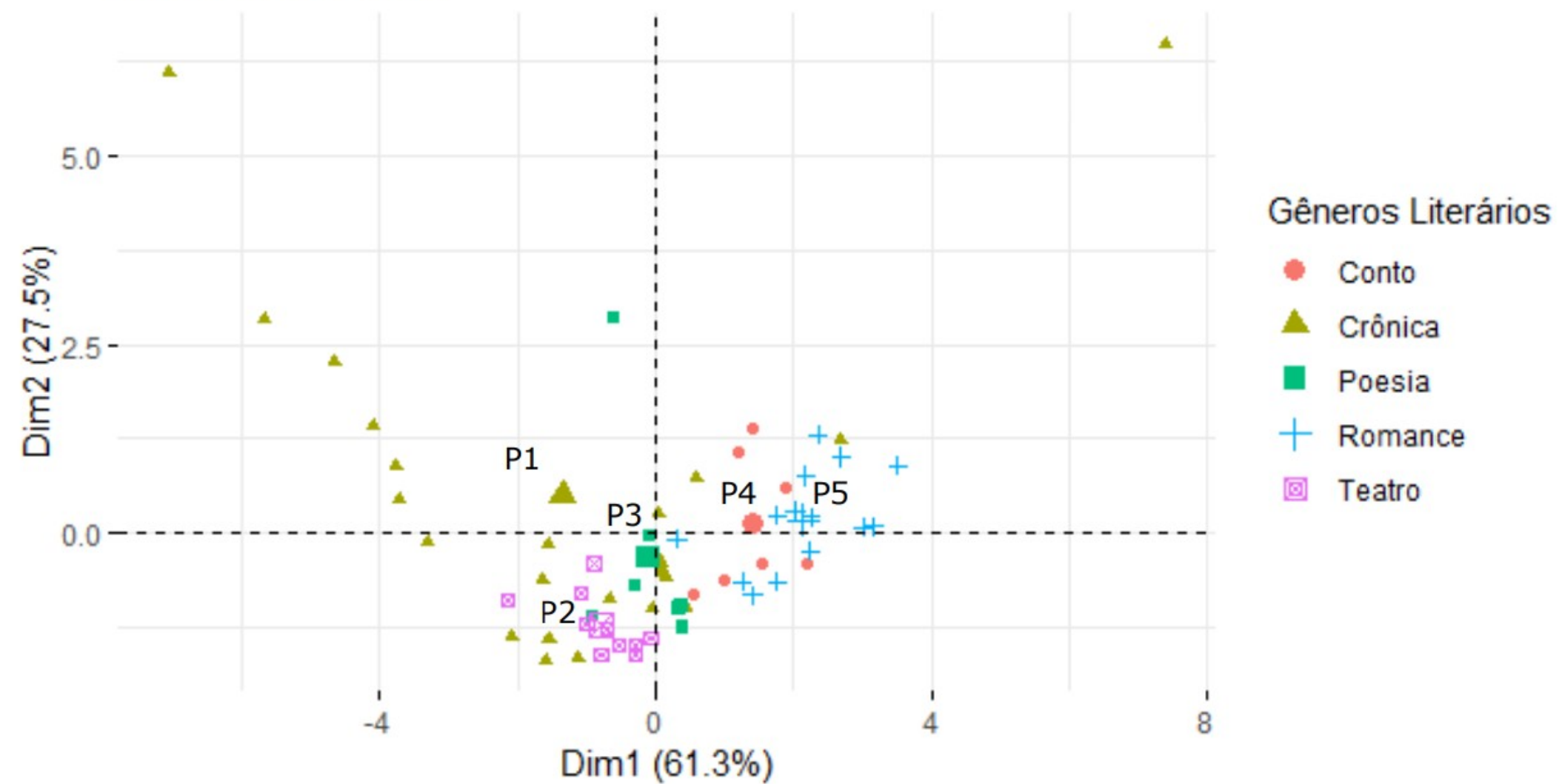


(b)

Janela de conexão 5



Individuals - PCA



Comentários finais...

Com a visualização das redes de co-ocorrência é possível realizar uma análise mais qualitativa das obras literárias, sendo evidenciadas palavras de maior relevância, o que torna mais fácil detectar conceitos relacionados com a temática de cada texto.

- A partir das métricas obtidas das redes de co-ocorrência, é possível realizar uma análise mais quantitativa, o que permitiu buscar padrões de similaridade e comparar autoria e gênero literário.
- A técnica utilizada foi a de análise de componentes principais sendo que, para ambas janelas de conexão, obtivemos uma porcentagem de explicação da variância de cerca de 90% com duas dimensões.