

# Topic modeling as a community-detection problem

Yuanming Tao<sup>1</sup>   Martin Gerlach<sup>2 4</sup>   Tiago P. Peixoto<sup>3 5</sup>   Eduardo G. Altmann<sup>1</sup>

<sup>1</sup>School of Mathematics and Statistics, University of Sydney, 2006 New South Wales, Australia.

<sup>2</sup>Department of Chemical and Biological Engineering, Northwestern University, Evanston, IL 60208, USA.

<sup>3</sup>Institute for Scientific Interchange Foundation, Via Alassio 11/c, 10126 Torino, Italy.

<sup>4</sup>Wikimedia Research Group, Berlin, Germany.

<sup>5</sup>Central European University, Budapest, Hungary Vienna, Austria.

NTU Singapore, 2019



CENTRE FOR COMPLEX SYSTEMS





## Topic models

- ▶ **Discover** what the books are talking about *automatically*



## Topic models

- ▶ **Discover** what the books are talking about *automatically*





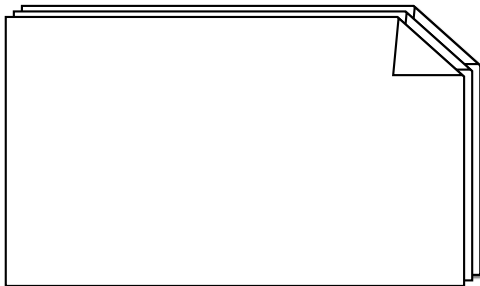
## Topic models

- ▶ **Discover** what the books are talking about *automatically*
- ▶ **Organize** the books onto bookshelves

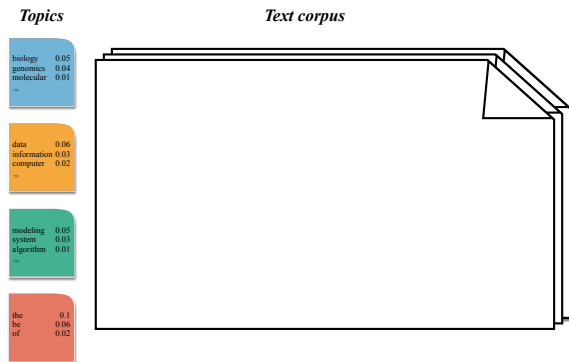
# TOPIC MODELING: A CONCRETE EXAMPLE

# TOPIC MODELING: A CONCRETE EXAMPLE

*Text corpus*



# TOPIC MODELING: A CONCRETE EXAMPLE



- Each **corpus-wide** topic is a cluster of words.

# TOPIC MODELING: A CONCRETE EXAMPLE

## Topics

biology	0.05
genomics	0.04
molecular	0.01
...	

data	0.06
information	0.03
computer	0.02
...	

modeling	0.05
system	0.03
algorithm	0.01
...	

the	0.1
be	0.06
of	0.02

## Text corpus

### Computational Biology

Computational biology involves the development and application of data-analytical and theoretical methods, mathematical modeling and computational simulation techniques to the study of biological, ecological, behavioral, and social systems.

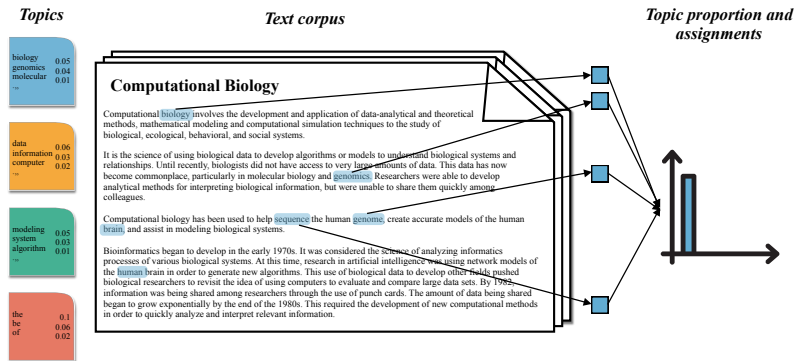
It is the science of using biological data to develop algorithms or models to understand biological systems and relationships. Until recently, biologists did not have access to very large amounts of data. This data has now become commonplace, particularly in molecular biology and genomics. Researchers were able to develop analytical methods for interpreting biological information, but were unable to share them quickly among colleagues.

Computational biology has been used to help sequence the human genome, create accurate models of the human brain, and assist in modeling biological systems.

Bioinformatics began to develop in the early 1970s. It was considered the science of analyzing informatics processes of various biological systems. At this time, research in artificial intelligence was using network models of the human brain in order to generate new algorithms. This use of biological data to develop other fields pushed biological researchers to revisit the idea of using computers to evaluate and compare large data sets. By 1982, information was being shared among researchers through the use of punch cards. The amount of data being shared began to grow exponentially by the end of the 1980s. This required the development of new computational methods in order to quickly analyze and interpret relevant information.

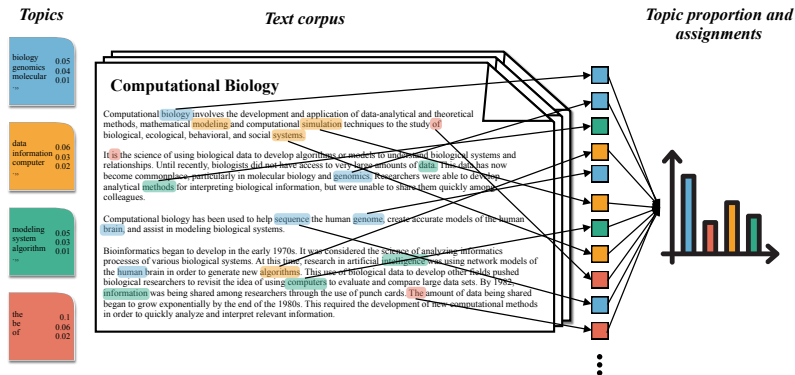
- Each **corpus-wide** topic is a cluster of words.

# TOPIC MODELING: A CONCRETE EXAMPLE



- ▶ Each **corpus-wide** topic is a cluster of words.
- ▶ Each document is a distribution over topics.

# TOPIC MODELING: A CONCRETE EXAMPLE



- ▶ Each **corpus-wide** topic is a cluster of words.
- ▶ Each document is a distribution over topics.

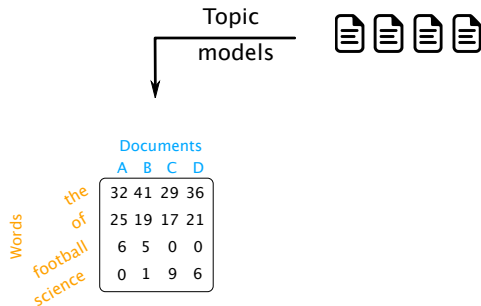
# LITERATURE REVIEW



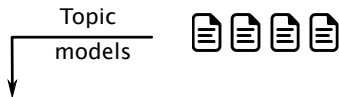
# LITERATURE REVIEW



# LITERATURE REVIEW



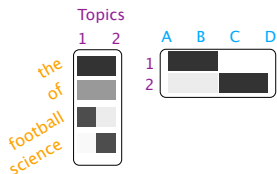
# LITERATURE REVIEW



Documents

	A	B	C	D
the	32	41	29	36
of	25	19	17	21
football	6	5	0	0
science	0	1	9	6

Words



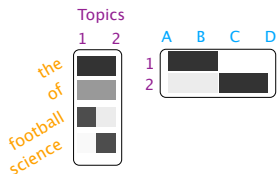
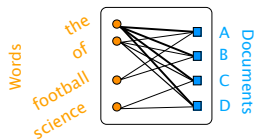
# LITERATURE REVIEW



Documents

	A	B	C	D
the	32	41	29	36
of	25	19	17	21
football	6	5	0	0
science	0	1	9	6

Words



# LITERATURE REVIEW

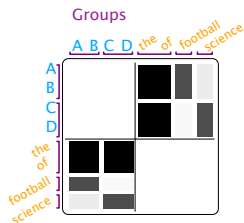
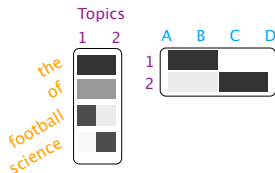
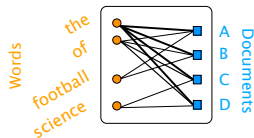
Topic  
models



Community  
detection

Documents

	A	B	C	D
the	32	41	29	36
of	25	19	17	21
football	6	5	0	0
science	0	1	9	6



Limitations: the above two models are *only* based on **word frequencies**.

# THE RESEARCH PROBLEM

Can we incorporate **further information** about documents?

# EXAMPLE: A WIKIPEDIA ARTICLE



# EXAMPLE: A WIKIPEDIA ARTICLE

## Computational biology

From Wikipedia, the free encyclopedia

**Computational biology** involves the development and application of data-analytical and theoretical methods, mathematical modeling and computational simulation techniques to the study of biological, ecological, behavioral, and social systems.<sup>[1]</sup> The field is broadly defined and includes foundations in [biology](#), [applied mathematics](#), [statistics](#), [biochemistry](#), [chemistry](#), [biophysics](#), [molecular biology](#), [genetics](#), [genomics](#), [computer science](#) and [evolution](#).<sup>[2]</sup>

Computational biology is different from [biological computing](#), which is a subfield of [computer science](#) and [computer engineering](#) using [bioengineering](#) and [biology](#) to build [computers](#), but is similar to [bioinformatics](#), which is an interdisciplinary science using computers to store and process biological data.

### Introduction

Computational Biology, which includes many aspects of [bioinformatics](#), is the science of using biological data to develop [algorithms](#) or [models](#) to understand biological systems and relationships. Until recently, biologists did not have access to very large amounts of data. This data has now become commonplace, particularly in [molecular biology](#) and [genomics](#). Researchers were able to develop analytical methods for interpreting biological information, but were unable to share them quickly among colleagues.<sup>[3]</sup>

Bioinformatics began to develop in the early 1970s. It was considered the science of analyzing informatics processes of various biological systems. At this time, research in [artificial intelligence](#) was using network models of the human brain in order to generate new [algorithms](#). This use of biological data to develop other fields pushed biological researchers to revisit the idea of using computers to evaluate and compare large data sets. By 1982, information was being shared among researchers through the use of punch cards. The amount of data being shared began to grow exponentially by the end of the 1980s. This required the development of new computational methods in order to quickly analyze and interpret relevant information.<sup>[3]</sup>

## EXAMPLE: A WIKIPEDIA ARTICLE

# Computational biology

 Connected to: [Genetics](#) [Molecular biology](#) [Chemistry](#)

From Wikipedia, the free encyclopedia

**Computational biology** involves the development and application of data-analytical and theoretical methods, mathematical modeling and computational simulation techniques to the study of biological, ecological, behavioral, and social systems.<sup>[1]</sup> The field is broadly defined and includes foundations in [biology](#), [applied mathematics](#), [statistics](#), [biochemistry](#), [chemistry](#), [biophysics](#), [molecular biology](#), [genetics](#), [genomics](#), [computer science](#) and [evolution](#).<sup>[2]</sup>

Computational biology is different from [biological computing](#), which is a subfield of [computer science](#) and [computer engineering](#) using [bioengineering](#) and [biology](#) to build [computers](#), but is similar to [bioinformatics](#), which is an interdisciplinary science using computers to store and process biological data.

## Introduction

Computational Biology, which includes many aspects of [bioinformatics](#), is the science of using biological data to develop [algorithms](#) or [models](#) to understand biological systems and relationships. Until recently, biologists did not have access to very large amounts of data. This data has now become commonplace, particularly in [molecular biology](#) and [genomics](#). Researchers were able to develop analytical methods for interpreting biological information, but were unable to share them quickly among colleagues.<sup>[3]</sup>

Bioinformatics began to develop in the early 1970s. It was considered the science of analyzing informatics processes of various biological systems. At this time, research in [artificial intelligence](#) was using network models of the human brain in order to generate new [algorithms](#). This use of biological data to develop other fields pushed biological researchers to revisit the idea of using computers to evaluate and compare large data sets. By 1982, information was being shared among researchers through the use of punch cards. The amount of data being shared began to grow exponentially by the end of the 1980s. This required the development of new computational methods in order to quickly analyze and interpret relevant information.<sup>[4]</sup>

Articles can have user-generated labels, called **metadata tags**.

# EXAMPLE: A WIKIPEDIA ARTICLE

The screenshot shows the top portion of a Wikipedia article. The title 'Computational biology' is at the top, followed by a 'Connected to' section with tags for 'Genetics', 'Molecular biology', and 'Chemistry'. Below this is the text 'From Wikipedia, the free encyclopedia'. The main text begins with a paragraph defining computational biology. A mouse cursor is hovering over the word 'Molecular biology' in the second paragraph. To the right of the text is a blue ribbon diagram of a protein structure. At the bottom of the screenshot, a 'Pro tip' is visible: 'Ctrl+Click any word on any website for quick definition'.

## Computational biology

Connected to: [Genetics](#) [Molecular biology](#) [Chemistry](#)

From Wikipedia, the free encyclopedia

**Computational biology** involves the development and application of data-analytical and theoretical methods, mathematical modeling and computational simulation techniques to the study of biological, ecological, behavioral, and social systems.<sup>[1]</sup> The field is broadly defined and includes foundations in [biology](#), [applied mathematics](#), [statistics](#), [biochemistry](#), [chemistry](#), [biophysics](#), [molecular biology](#), [genetics](#), [genomics](#), [computer science](#) and [evolution](#).<sup>[1]</sup>

Computational biology is an interdisciplinary field of [computer science](#) and [computer engineering](#) applied to [biology](#) or [bioinformatics](#), which is an interdisciplinary field of [computer science](#) and [informatics](#) that uses computational methods to analyze biological data.

### Introduction

Computational biology is a branch of biology that concerns the molecular basis of biological activity between biomolecules in the various systems of a cell, including the interactions between DNA, RNA, proteins and their biosynthesis, as well as the regulation of these interactions. Writing in *Nature* in 1961, William Astbury described molecular biology as: "...not so much a technique as an approach, an approach from the viewpoint of the so-called basic sciences with the leading idea of searching below the large-scale manifestations of classical biology for the corresponding molecular plan. It is concerned particularly with the forms of biological molecules and [...] is predominantly three-dimensional and structural – which does not mean, however, that it is merely a refinement of morphology.

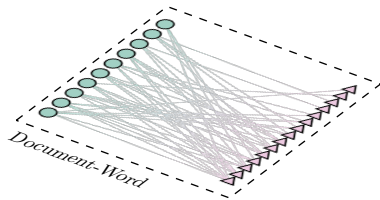
Pro tip: Ctrl+Click any word on any website for quick definition

Articles can contain external links to other articles, e.g. **hyperlinks**.

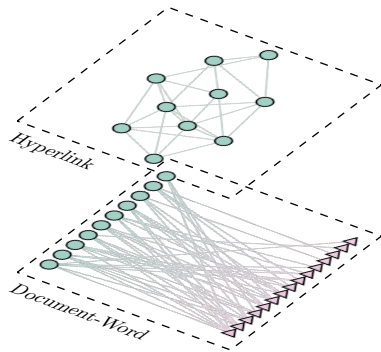
How do we incorporate **metadata tags** and **hyperlinks** in the model?

# MODELING TEXTS WITH AUXILIARY INFORMATION USING MULTILAYERED NETWORK

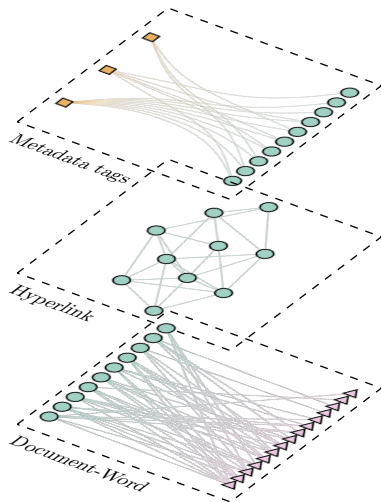
# MODELING TEXTS WITH AUXILIARY INFORMATION USING MULTILAYERED NETWORK



# MODELING TEXTS WITH AUXILIARY INFORMATION USING MULTILAYERED NETWORK



# MODELING TEXTS WITH AUXILIARY INFORMATION USING MULTILAYERED NETWORK





What is the model for the formation of the above network structures?

# THE STOCHASTIC BLOCK MODELS (SBM)

# THE STOCHASTIC BLOCK MODELS (SBM)

**Planted partition:**  $N$  nodes divided into  $B$  blocks.

# THE STOCHASTIC BLOCK MODELS (SBM)

**Planted partition:**  $N$  nodes divided into  $B$  blocks.

Parameters:

$b_i \rightarrow$  block membership of node  $i$

$p_{rs} \rightarrow$  probability of an edge between nodes of groups  $r$  and  $s$

# THE STOCHASTIC BLOCK MODELS (SBM)

**Planted partition:**  $N$  nodes divided into  $B$  blocks.

Parameters:

$b_i \rightarrow$  block membership of node  $i$

$p_{rs} \rightarrow$  probability of an edge between nodes of groups  $r$  and  $s$

# THE STOCHASTIC BLOCK MODELS (SBM)

**Planted partition:**  $N$  nodes divided into  $B$  blocks.

Parameters:

$b_i \rightarrow$  block membership of node  $i$

$p_{rs} \rightarrow$  probability of an edge between nodes of groups  $r$  and  $s$

Likelihood:  $P(A_{ij}|b_i, p_{b_i, b_j}) = \text{Poisson}(C \times p_{b_i, b_j})$ , where  $C$  is a constant.

# THE STOCHASTIC BLOCK MODELS (SBM)

**Planted partition:**  $N$  nodes divided into  $B$  blocks.

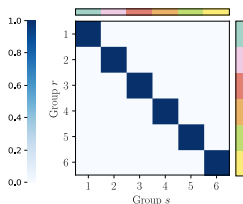
Parameters:

$b_i \rightarrow$  block membership of node  $i$

$p_{rs} \rightarrow$  probability of an edge between nodes of groups  $r$  and  $s$

Likelihood:  $P(A_{ij}|b_i, p_{b_i, b_j}) = \text{Poisson}(C \times p_{b_i, b_j})$ , where  $C$  is a constant.

**Example:**



# THE STOCHASTIC BLOCK MODELS (SBM)

**Planted partition:**  $N$  nodes divided into  $B$  blocks.

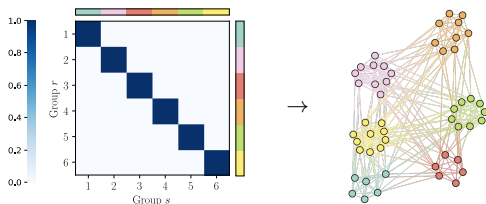
Parameters:

$b_i \rightarrow$  block membership of node  $i$

$p_{rs} \rightarrow$  probability of an edge between nodes of groups  $r$  and  $s$

Likelihood:  $P(A_{ij}|b_i, p_{b_i, b_j}) = \text{Poisson}(C \times p_{b_i, b_j})$ , where  $C$  is a constant.

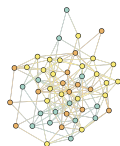
**Example:**



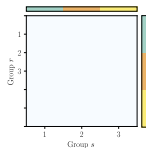


# THE SBM IS NOT RESTRICTED TO ASSORTATIVE STRUCTURES

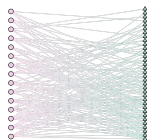
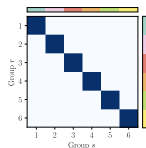
# THE SBM IS NOT RESTRICTED TO ASSORTATIVE STRUCTURES



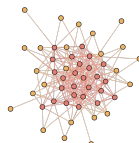
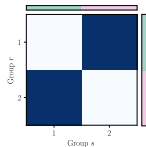
(a) The random network



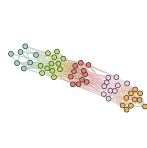
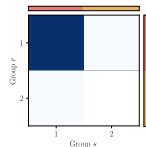
(b) The assortative structure



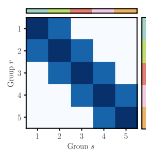
(c) The bipartite structure



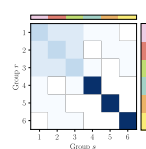
(d) The core-periphery structure



(e) The ordered structure



(f) The mixed pattern



# SBM WITH INDEPENDENT LAYERS

# SBM WITH INDEPENDENT LAYERS

## Parameter:

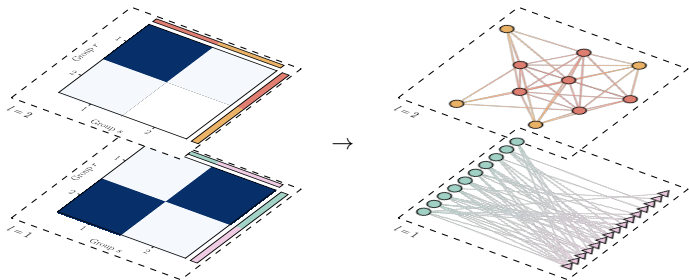
$z_{il} \rightarrow$  binary values to indicate whether a node  $i$  appears on layer  $l$

# SBM WITH INDEPENDENT LAYERS

## Parameter:

$z_{il} \rightarrow$  binary values to indicate whether a node  $i$  appears on layer  $l$

## Example:

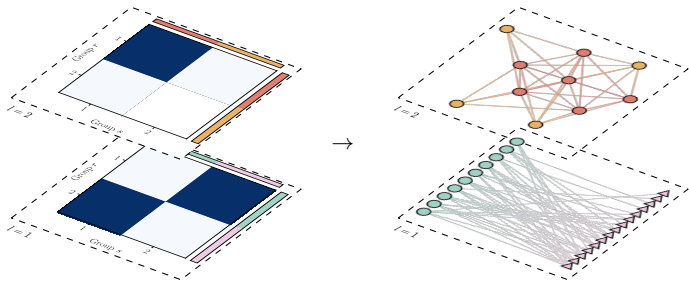


# SBM WITH INDEPENDENT LAYERS

## Parameter:

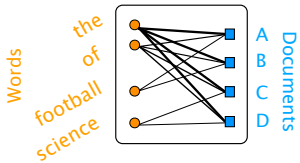
$z_{il} \rightarrow$  binary values to indicate whether a node  $i$  appears on layer  $l$

## Example:

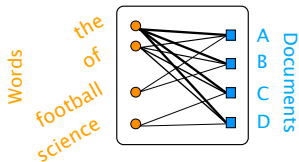


## Important assumptions:

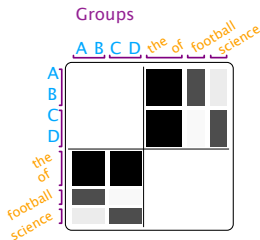
Degree distribution varies across different layers.



# THE BAYESIAN INFERENCE OF THE SBM

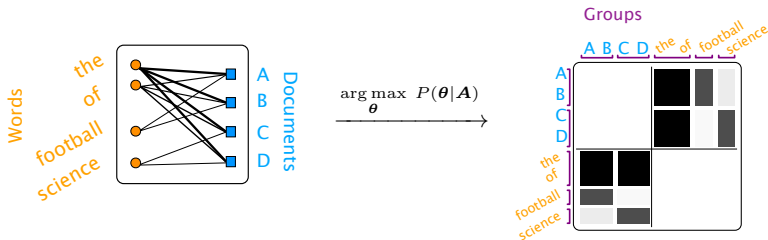


$$\arg \max_{\theta} P(\theta | A)$$





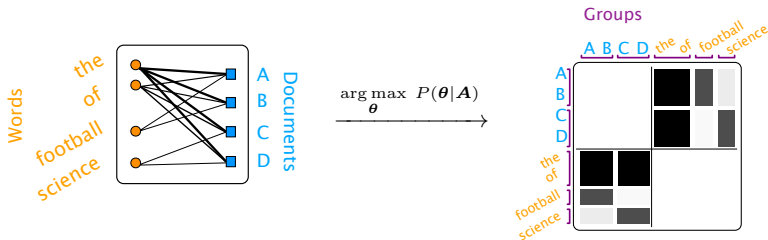
# THE BAYESIAN INFERENCE OF THE SBM



$$\Sigma = -\ln P(\theta|A)$$

$\Sigma \rightarrow$  Total information necessary to describe the data.

# THE BAYESIAN INFERENCE OF THE SBM



$$\Sigma = -\ln P(\theta|\mathbf{A}) = \underbrace{-\ln P(\mathbf{A}|\theta)}_{\text{data|model, } \mathcal{S}} - \underbrace{\ln P(\theta)}_{\text{model, } \mathcal{L}}$$

$\Sigma \rightarrow$  Total information necessary to describe the data.

$\mathcal{S} \rightarrow$  Information required to describe the network  $\mathbf{A}$ , when the model is known.

$\mathcal{L} \rightarrow$  Information required to describe the model parameters  $\theta$ .

Will incorporating additional information available about the documents help us to **better** classify them?

# CASE STUDY: THE WIKIPEDIA ARTICLES

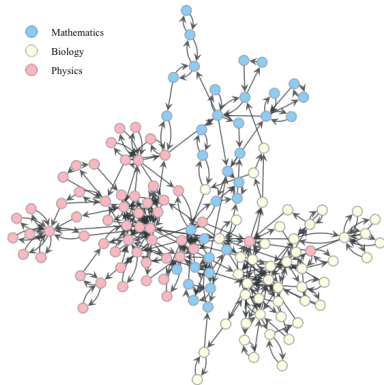
# CASE STUDY: THE WIKIPEDIA ARTICLES

- ▶ 138 Wikipedia articles with three categories:

mathematics: 34,

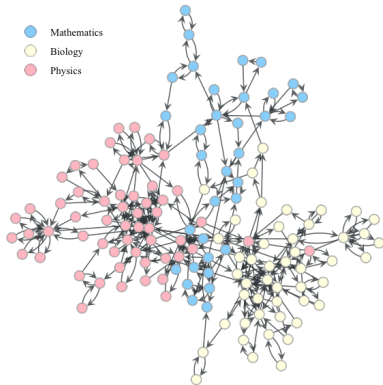
physics: 56,

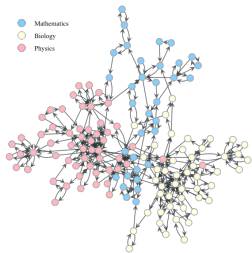
biology: 48

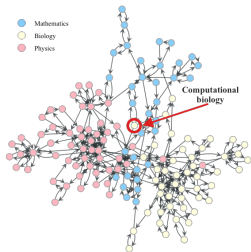


# CASE STUDY: THE WIKIPEDIA ARTICLES

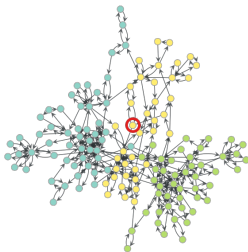
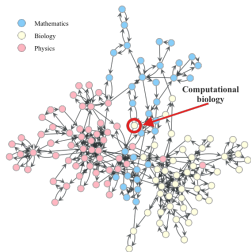
- ▶ 138 Wikipedia articles with three categories:
  - mathematics: 34,
  - physics: 56,
  - biology: 48
- ▶ Number of hyperlinks: 341
- ▶ Number of distinct words: 16,378
- ▶ Number of words: 351,710

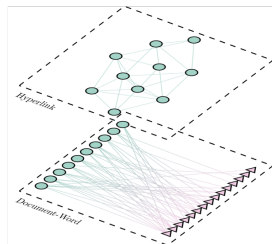
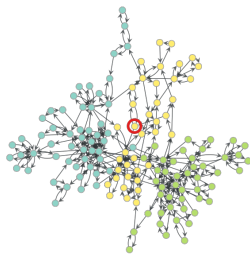
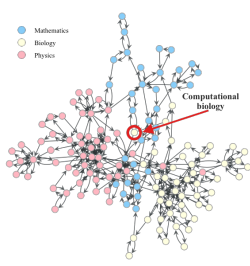


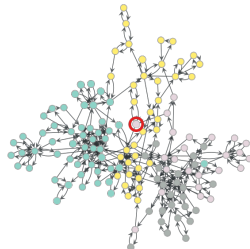
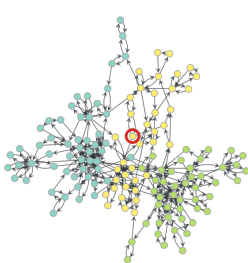
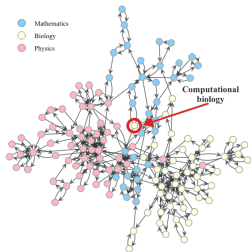




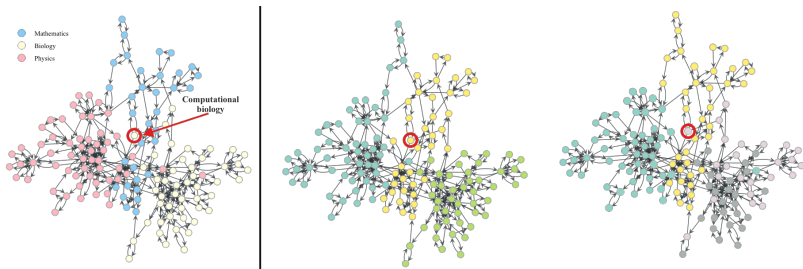




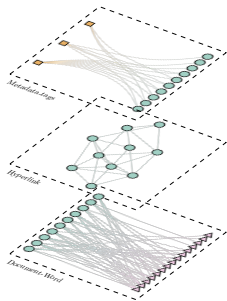
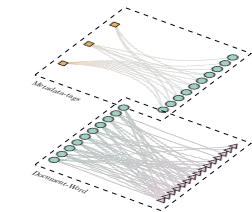
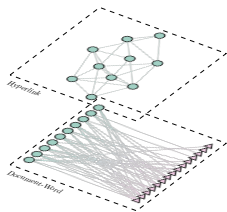
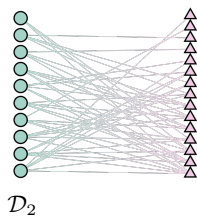
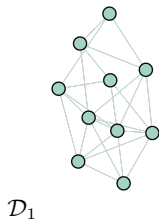
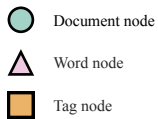




# MEASURING PARTITION SIMILARITY: NORMALIZED MUTUAL INFORMATION (NMI)



- ▶  $NMI \in [0, 1]$ .
- ▶ Larger NMI values indicate better agreement.



In each dataset  $\mathcal{D}_1, \dots, \mathcal{D}_5$

- ▶ Run the MCMC algorithm for multiple times.
- ▶ Obtain the inferred partition after annealing.

# COMPARISON OF PARTITION SIMILARITIES OF DOCUMENTS

Wikipedia Labels	1
$\mathcal{D}_1$	0.51

In each dataset  $\mathcal{D}_1, \dots, \mathcal{D}_5$

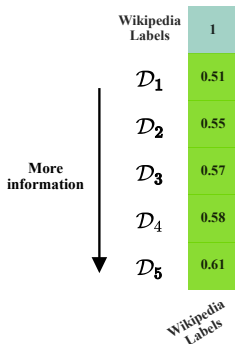
- ▶ Run the MCMC algorithm for multiple times.
- ▶ Obtain the inferred partition after annealing.

$\mathcal{D}_1$ : **without** word nodes.

# COMPARISON OF PARTITION SIMILARITIES OF DOCUMENTS

In each dataset  $\mathcal{D}_1, \dots, \mathcal{D}_5$

- ▶ Run the MCMC algorithm for multiple times.
- ▶ Obtain the inferred partition after annealing.



$\mathcal{D}_1$ : **without** word nodes.

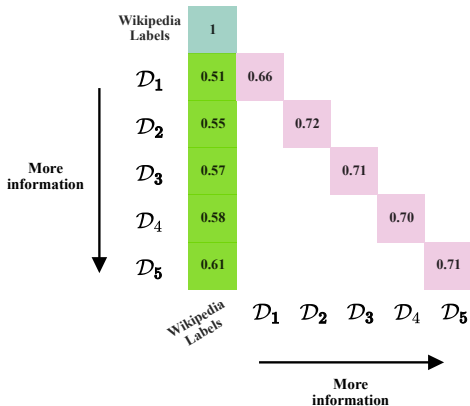
$\mathcal{D}_2 \sim \mathcal{D}_5$ : **with** word nodes.



# COMPARISON OF PARTITION SIMILARITIES OF DOCUMENTS

In each dataset  $\mathcal{D}_1, \dots, \mathcal{D}_5$

- ▶ Run the MCMC algorithm for multiple times.
- ▶ Obtain the inferred partition after annealing.



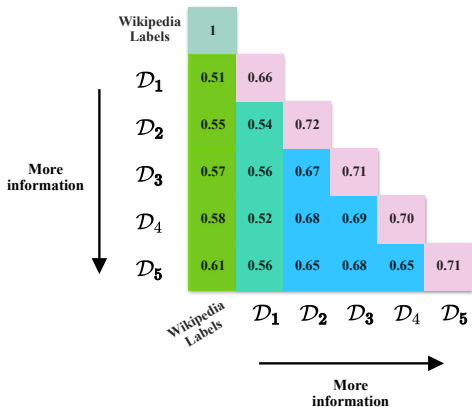
$\mathcal{D}_1$ : **without** word nodes.

$\mathcal{D}_2 \sim \mathcal{D}_5$ : **with** word nodes.

# COMPARISON OF PARTITION SIMILARITIES OF DOCUMENTS

In each dataset  $\mathcal{D}_1, \dots, \mathcal{D}_5$

- ▶ Run the MCMC algorithm for multiple times.
- ▶ Obtain the inferred partition after annealing.



$\mathcal{D}_1$ : **without** word nodes.

$\mathcal{D}_2 \sim \mathcal{D}_5$ : **with** word nodes.

Main message:

## Main message:

- ▶ We extend the previous work by incorporating additional information available about documents as additional layers in the same SBM framework.

## Main message:

- ▶ We extend the previous work by incorporating additional information available about documents as additional layers in the same SBM framework.
- ▶ Compared to Wikipedia labeling of articles, incorporating more information will lead to better agreement.

## Main message:

- ▶ We extend the previous work by incorporating additional information available about documents as additional layers in the same SBM framework.
- ▶ Compared to Wikipedia labeling of articles, incorporating more information will lead to better agreement.
- ▶ Document-word layer is dominating in the inference.

## Main message:

- ▶ We extend the previous work by incorporating additional information available about documents as additional layers in the same SBM framework.
- ▶ Compared to Wikipedia labeling of articles, incorporating more information will lead to better agreement.
- ▶ Document-word layer is dominating in the inference.
- ▶ Partitions of documents are more similar when incorporating words.

## Main message:

- ▶ We extend the previous work by incorporating additional information available about documents as additional layers in the same SBM framework.
- ▶ Compared to Wikipedia labeling of articles, incorporating more information will lead to better agreement.
- ▶ Document-word layer is dominating in the inference.
- ▶ Partitions of documents are more similar when incorporating words.

## References:

1. Tiago P. Peixoto. "Bayesian stochastic blockmodeling." arXiv preprint arXiv:1705.10225 (2017).
2. Martin Gerlach, Tiago P. Peixoto, and Eduardo G. Altmann. "A network approach to topic models." Science advances 4, no. 7 (2018): eaaq1360.
3. The graph-tool package <https://graph-tool.skewed.de>.