

Statistical Learning with Networks and Texts

Charles BOUVEYRON

Professor of Statistics
Chair of Excellence Inria on "Data Science"

Laboratoire LJAD, UMR CNRS 7351
Equipe Asclepios, Inria Sophia-Antipolis
Université Côte d'Azur

charles.bouveyron@unice.fr
@cbouveyron



Preamble

“Essentially, all models are wrong but some are useful”

George E.P. Box

Outline

1. Introduction
2. The Stochastic Topic Block Model
3. Numerical application: The Enron case
4. The Linkage project
5. Conclusion

Introduction

In statistical learning, the challenge nowadays is to learn from data which are:

- high-dimensional (p large),
- big or as stream (n large),
- evolutive (evolving phenomenon),
- heterogeneous (categorical, functional, networks, texts, ...)

Introduction

In statistical learning, the challenge nowadays is to learn from data which are:

- high-dimensional (p large),
- big or as stream (n large),
- evolutive (evolving phenomenon),
- heterogeneous (categorical, functional, networks, texts, ...)

In any case, the understanding of the results is essential :

- the practitioners are interested in **visualizing** or **clustering** their data,
- to have a selection of the **relevant original variables** for interpretation,
- and to have a **probabilistic model** supposed to have generated the data.

Introduction

Statistical analysis of (social) networks has become a strong discipline:

- description and comparison of networks,
- network visualization,
- clustering of network nodes.

Introduction

Statistical analysis of (social) networks has become a strong discipline:

- description and comparison of networks,
- network visualization,
- clustering of network nodes.

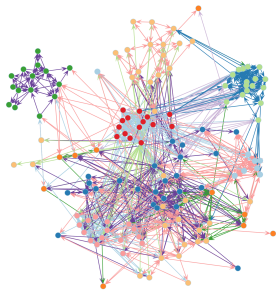
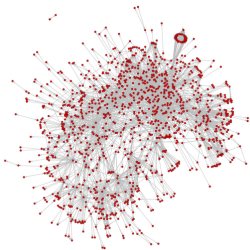
with applications in domains ranging from biology to historical sciences:

- biology: analysis of gene regulation processes,
- social sciences: analysis of political blogs,
- historical sciences: clustering and comparison of medieval social networks
 - Bouveyron, Lamassé et al., *The random subgraph model for the analysis of an ecclesiastical network in merovingian Gaul*, *The Annals of Applied Statistics*, vol. 8(1), pp. 377-405, 2014.

Introduction

Networks can be observed **directly** or **indirectly** from a variety of sources:

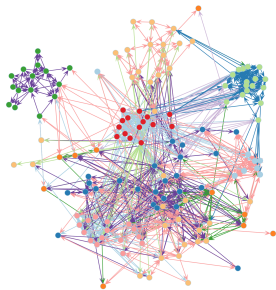
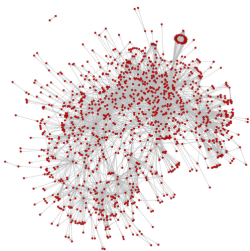
- social websites (Facebook, Twitter, ...),
- personal emails (from your Gmail, Clinton's mails, ...),
- emails of a company (Enron Email data),
- digital/numeric documents (Panama papers, co-authorships, ...),
- and even archived documents in libraries (digital humanities).



Introduction

Networks can be observed **directly** or **indirectly** from a variety of sources:

- social websites (Facebook, Twitter, ...),
- personal emails (from your Gmail, Clinton's mails, ...),
- emails of a company (Enron Email data),
- digital/numeric documents (Panama papers, co-authorships, ...),
- and even archived documents in libraries (digital humanities).



⇒ most of these sources involve text!

An introductory example

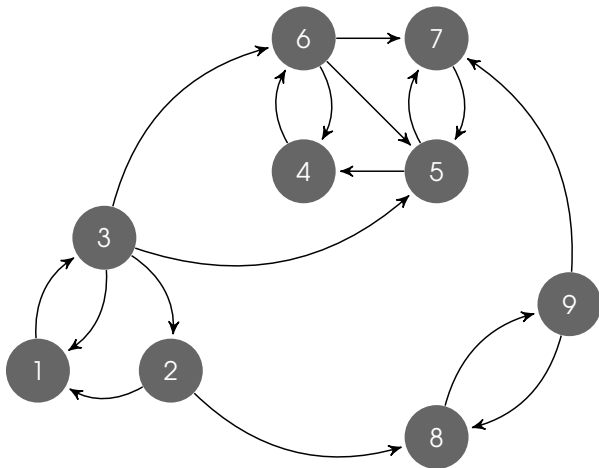


Figure: An (hypothetic) email network between a few individuals.

An introductory example

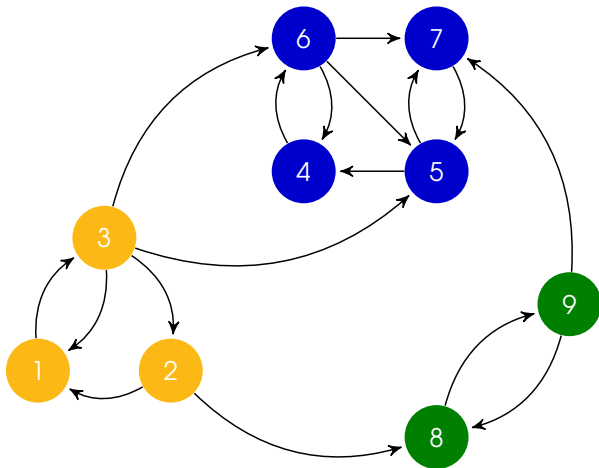


Figure: A typical clustering result for the (directed) binary network.

An introductory example

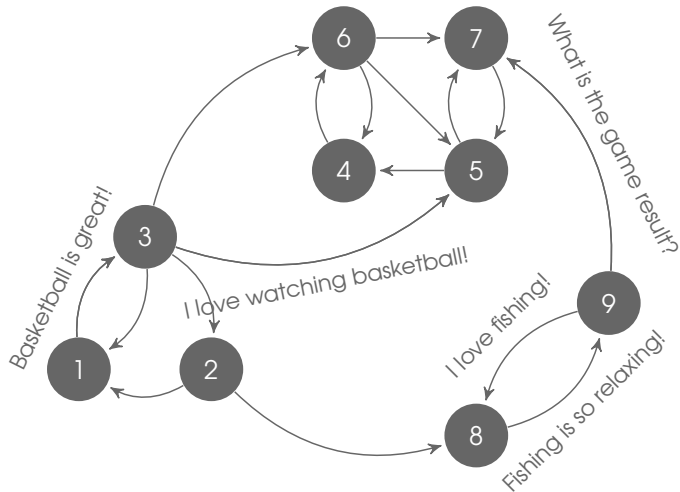


Figure: The (directed) network with textual edges.

An introductory example

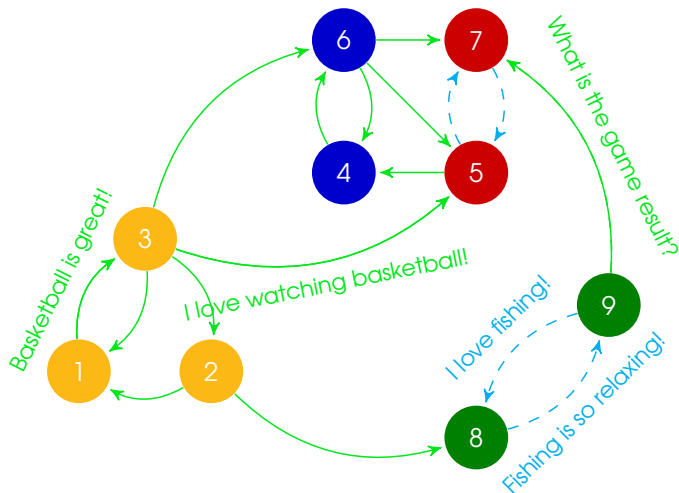


Figure: Expected clustering result for the (directed) network with textual edges.

Outline

1. Introduction
2. The Stochastic Topic Block Model
3. Numerical application: The Enron case
4. The Linkage project
5. Conclusion

STBM : Context and notations

We are interesting in **clustering the nodes** of a (directed) network of M vertices into Q groups:

- the network is represented by its $M \times M$ **adjacency matrix** A :

$$A_{ij} = \begin{cases} 1 & \text{if there is an edge between } i \text{ and } j \\ 0 & \text{otherwise} \end{cases}$$

- if $A_{ij} = 1$, the textual edge is characterized by a set of D_{ij} **documents**:

$$W_{ij} = (W_{ij}^1, \dots, W_{ij}^d, \dots, W_{ij}^{D_{ij}}),$$

- each document W_{ij}^d is made of N_{ij}^d **words**:

$$W_{ij}^d = (W_{ij}^{d1}, \dots, W_{ij}^{dn}, \dots, W_{ij}^{dN_{ij}^d}).$$

STBM : Modeling of the edges

Let us assume that edges are generated according to a SBM model:

- each node i is associated with an (unobserved) group among Q according to:

$$Y_i \sim \mathcal{M}(\rho),$$

where $\rho \in [0, 1]^Q$ is the vector of group proportions,

STBM : Modeling of the edges

Let us assume that edges are generated according to a SBM model:

- each node i is associated with an (unobserved) group among Q according to:

$$Y_i \sim \mathcal{M}(\rho),$$

where $\rho \in [0, 1]^Q$ is the vector of group proportions,

- the presence of an edge A_{ij} between i and j is drawn according to:

$$A_{ij} | Y_{iq} Y_{jr} = 1 \sim \mathcal{B}(\pi_{qr}),$$

where $\pi_{qr} \in [0, 1]$ is the connection probability between clusters q and r .

STBM : Modeling of the documents

The generative model for the documents is as follows:

- each pair of clusters (q, r) is first associated to a **vector of topic proportions** $\theta_{qr} = (\theta_{qrk})_k$ sampled from a Dirichlet distribution:

$$\theta_{qr} \sim \text{Dir}(\alpha),$$

such that $\sum_{k=1}^K \theta_{qrk} = 1, \forall (q, r)$.

STBM : Modeling of the documents

The generative model for the documents is as follows:

- each pair of clusters (q, r) is first associated to a **vector of topic proportions** $\theta_{qr} = (\theta_{qrk})_k$ sampled from a Dirichlet distribution:

$$\theta_{qr} \sim \text{Dir}(\alpha),$$

such that $\sum_{k=1}^K \theta_{qrk} = 1, \forall (q, r)$.

- the n th word W_{ij}^{dn} of documents d in W_{ij} is then associated to a **latent topic vector** Z_{ij}^{dn} according to:

$$Z_{ij}^{dn} | \{A_{ij} Y_{iq} Y_{jr} = 1, \theta\} \sim \mathcal{M}(1, \theta_{qr}).$$

STBM : Modeling of the documents

The generative model for the documents is as follows:

- each pair of clusters (q, r) is first associated to a **vector of topic proportions** $\theta_{qr} = (\theta_{qrk})_k$ sampled from a Dirichlet distribution:

$$\theta_{qr} \sim \text{Dir}(\alpha),$$

such that $\sum_{k=1}^K \theta_{qrk} = 1, \forall (q, r)$.

- the n th word W_{ij}^{dn} of documents d in W_{ij} is then associated to a **latent topic vector** Z_{ij}^{dn} according to:

$$Z_{ij}^{dn} | \{A_{ij} Y_{iq} Y_{jr} = 1, \theta\} \sim \mathcal{M}(1, \theta_{qr}).$$

- then, given Z_{ij}^{dn} , the **word** W_{ij}^{dn} is assumed to be drawn from a multinomial distribution:

$$W_{ij}^{dn} | Z_{ij}^{dnk} = 1 \sim \mathcal{M}(1, \beta_k = (\beta_{k1}, \dots, \beta_{kV})),$$

where V is the vocabulary size.

STBM at a glance...

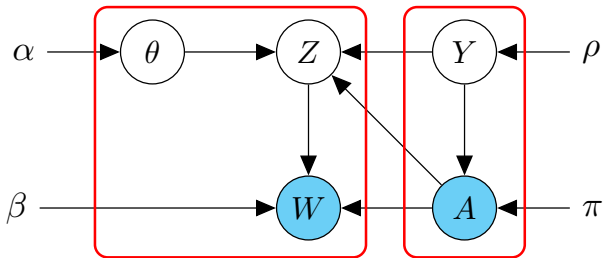


Figure: The stochastic topic block model.

The C-VEM algorithm for inference

The C-VEM algorithm is as follows:

- we use a VEM algorithm to maximize $\tilde{\mathcal{L}}$ with respect to β and $R(Z, \theta)$, which *essentially* corresponds to the VEM algorithm of Blei et al. (2003),
- then, $\log p(A, Y | \rho, \pi)$ is maximized with respect to ρ and π to provide estimates,
- finally, $\mathcal{L}(R(\cdot); Y, \rho, \pi, \beta)$ is maximized with respect to Y , which is the only term involved in both $\tilde{\mathcal{L}}$ and the SBM complete data log-likelihood.

The C-VEM algorithm for inference

The C-VEM algorithm is as follows:

- we use a VEM algorithm to maximize $\tilde{\mathcal{L}}$ with respect to β and $R(Z, \theta)$, which *essentially* corresponds to the VEM algorithm of Blei et al. (2003),
- then, $\log p(A, Y | \rho, \pi)$ is maximized with respect to ρ and π to provide estimates,
- finally, $\mathcal{L}(R(\cdot); Y, \rho, \pi, \beta)$ is maximized with respect to Y , which is the only term involved in both $\tilde{\mathcal{L}}$ and the SBM complete data log-likelihood.

Optimization over Y :

- we propose an online approach which cycles randomly through the vertices,
- at each step, a single vertex i is considered and all membership vectors $Y_{j \neq i}$ are held fixed,
- for vertex i , we look for every possible cluster assignment Y_i and the one which maximizes $\mathcal{L}(R(\cdot); Y, \rho, \pi, \beta)$ is kept.

Outline

1. Introduction
2. The Stochastic Topic Block Model
3. Numerical application: The Enron case
4. The Linkage project
5. Conclusion

Analysis of the Enron Emails

The Enron data set:

- all emails between 149 Enron employees,
- from 1999 to the bankrupt in late 2001,
- almost 253 000 emails in the whole data base.

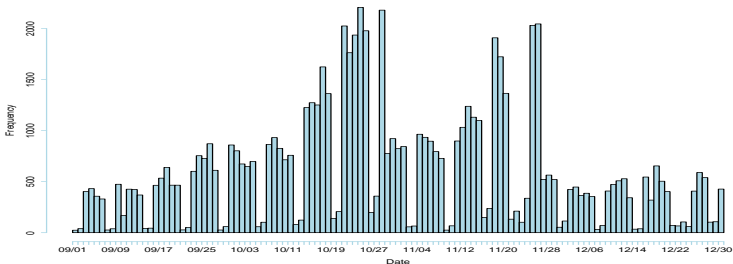


Figure: Temporal distribution of Enron emails.

Analysis of the Enron Emails

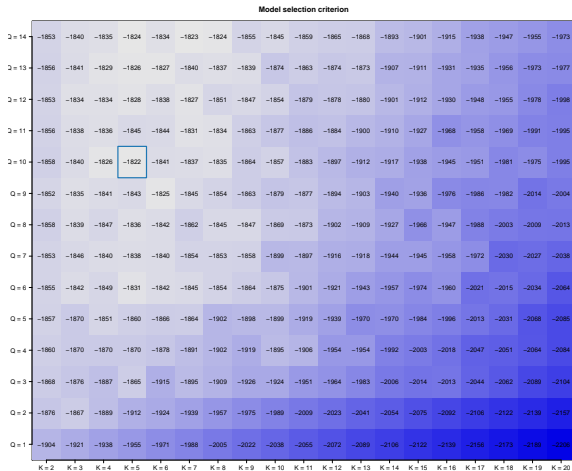


Figure: Model selection for STBM on the Enron network.

Analysis of the Enron Emails

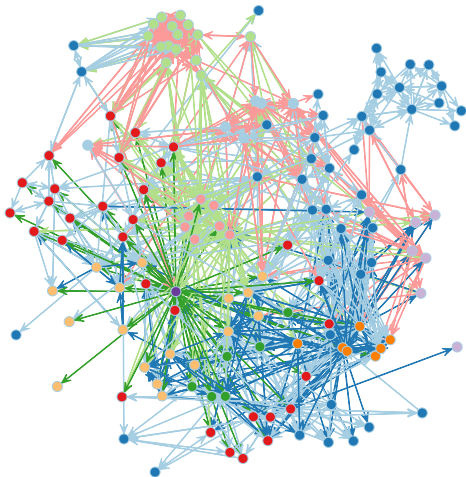


Figure: Clustering of the Enron network.

Analysis of the Enron Emails

cycle	grigsby	edison	backup	mmbtud
oto	afghanistan	puc	seat	harris
usage	viewing	interview	test	watson
select	desk	state	location	capacity
prorata	phillip	interviewers	building	transwestern
storage	ground	dwr	supplies	deliveries
interruptible	park	davis	computer	hayslett
declared	taleban	fantastic	announcement	master
equal	forces	dinner	notified	lynn
ridge	named	said	phones	socalgas
forecast	sheppard	saturday	seats	shackleton
windows	fundamental	super	locations	donoho
wheeler	tori	california	regular	lindy
nom	allen	mara	assignments	kay
injections	ermis	dasovich	rely	sara
elapsed	ina	phase	assignment	geaccone
limits	kuykendall	governor	numbers	pkgs
gas	gaskill	steffes	equipment	juan
receipt	bin	rto	aside	kilmer
clock	heizenrader	contracts	floors	netting
Topic 1	Topic 2	Topic 3	Topic 4	Topic 5

Figure: Most specific terms in the found topics for the Enron data.

Analysis of the Enron Emails

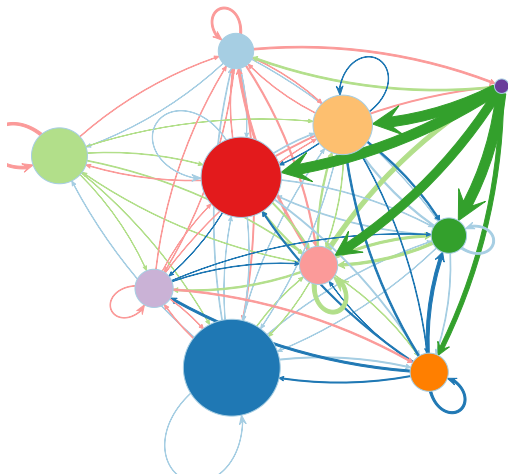


Figure: Meta-network for the Enron data set.

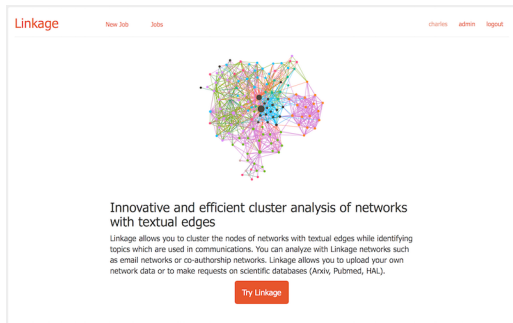
Outline

1. Introduction
2. The Stochastic Topic Block Model
3. Numerical application: The Enron case
4. The Linkage project
5. Conclusion

Innovation: the linkage project

From research to Innovation:

- the project is supported by SATT IDFIInnov,
- 50 k€ for SaaS platform www.linkage.fr
- 200k€ for further work (dynamic, sparsity)



Linkage

New Job Jobs

charles admin logout

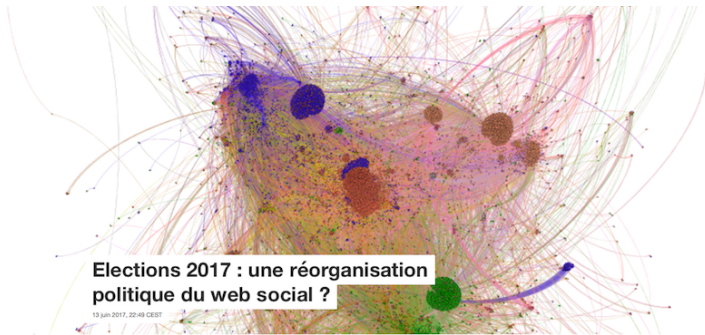
Innovative and efficient cluster analysis of networks with textual edges

Linkage allows you to cluster the nodes of networks with textual edges while identifying topics which are used in communications. You can analyze with Linkage networks such as email networks or co-authorship networks. Linkage allows you to upload your own network data or to make requests on scientific databases (Arxiv, Pubmed, HAL).

Try Linkage

www.linkage.fr

Analysis of the 2017 French presidential election



Réseaux des tweets des français liés à la politique des 17 et 18 avril. P. Latouche, CC BY

Adresse électronique

Twitter 13

Facebook 48

LinkedIn 11

Imprimer

Emmanuel Macron vient d'être élu à la présidence de la République sur un programme dont une des priorités est la recomposition de la vie politique. La période que nous traversons, entre les deux tours des législatives, est donc sujette à de fortes interrogations quant à la réorganisation à venir des partis politiques.

Afin d'apporter un éclairage sur ce point, nous avons étudié pendant les semaines qui ont précédé le second tour de l'élection présidentielle les mouvements et transferts entre les partis, avec un prisme particulier.

Auteurs



Pierre Latouche

Maître de conférences en Mathématiques Appliquées, Université Paris 1 Panthéon-Sorbonne



Charles Bouveyron

Professeur des Universités en Mathématiques Appliquées, Université Paris Descartes - USPC

up5.fr/presid2017

Outline

1. Introduction
2. The Stochastic Topic Block Model
3. Numerical application: The Enron case
4. The Linkage project
5. Conclusion

Conclusion

We proposed a **new statistical model, called STBM**, for :

- the clustering of the nodes of networks with textual edges,
- which also "clusters" the messages into general topics,
- it provides an effective summary of the whole data (network + texts).

STBM can be applied to:

- communication networks (emails, web forums, twitters, ...),
- co-authorship networks (scientific publications, patents, ...),
- and can even applied to networks with images (Instagram, ...).

Reference:

C. Bouveyron, P. Latouche and R. Zreik, The Stochastic Topic Block Model for the Clustering of Networks with Textual Edges, Statistics & Computing, in press, 2017.