



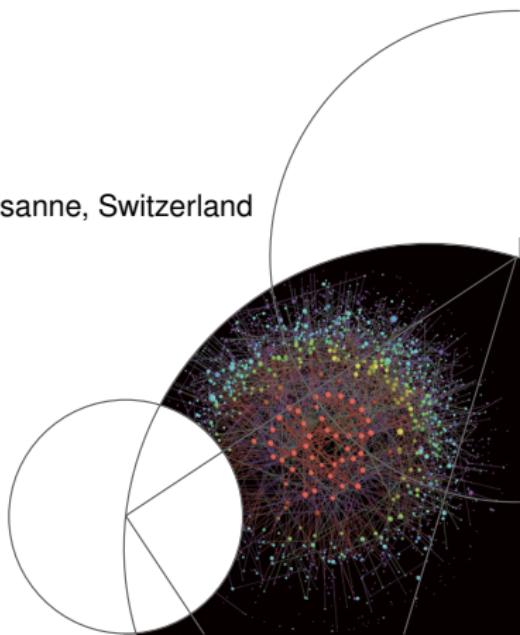
Automatic identification of relevant concepts in scientific publications

Alessio Cardillo

Laboratory for Statistical Biophysics (LBS)

École polytechnique fédérale de Lausanne (EPFL), Lausanne, Switzerland

<http://bifi.es/~cardillo/>



Acknowledgements

Main Collaborators

- Paolo De Los Rios (EPFL)
- Andrea Martini (EPFL)

Other Collaborators

- Alex Constantin (EPFL)
- Vasyl Palchykov, Valerio Gemmetto, Diego Garlaschelli (Leiden - The Netherlands)
- ScienceWISE Team (Lausanne)

JAEWm

newsblog

Nature brings you breaking news from the world of science

News & Comment > News blog Archive > Post

Previous post

Climate change is present danger, US warns

Next post

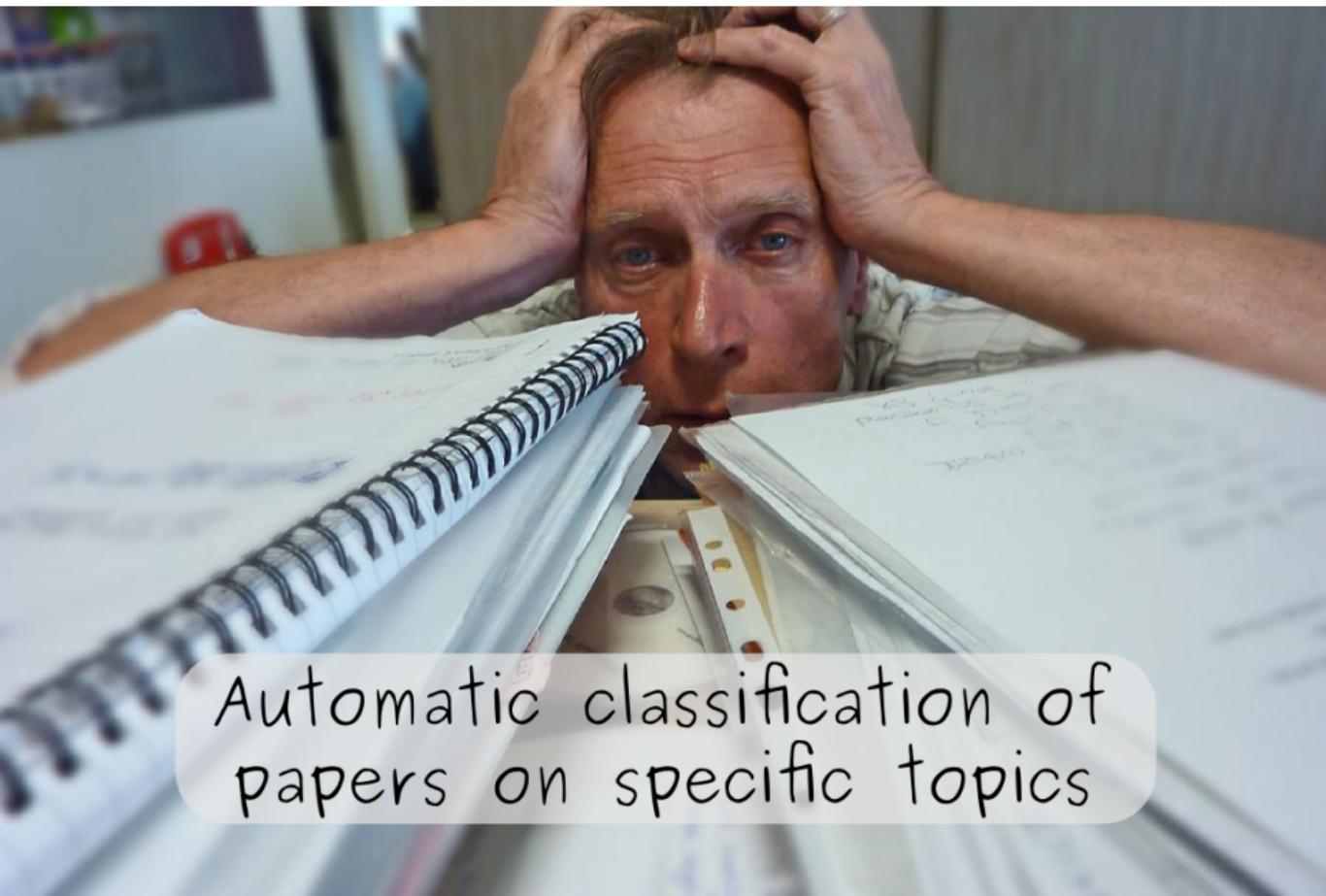
German research agencies condemn animal-rights attack on neuroscientist

NEWS BLOG

Global scientific output doubles every nine years

07 May 2014 | 16:46 GMT | Posted by Richard Van Noorden | Category: Policy, Publishing

It's a common complaint among academics: today's researchers are publishing too much, too fast. But just



Automatic classification of
papers on specific topics



Recent ontology graph



Recently bookmarked papers

Properties of a possible class of particles
astro-ph/9505117 Luis Gonzalez-Mestres

The apparent Lorentz invariance of the laws of physi...

Introduction to the Standard Model and E
0901.0241 Paul Langacker

A concise introduction is given to the standard model. Including the structure of the QCD and electroweak Lagrangians, spontaneous symmetry breaking, experimental tests, and problems.

[Standard Model](#) [Quantum chromodynamics](#)
[Weak interaction](#) ...

<http://sciencewise.info>

Outline

- Introduction on similarity networks.
- ★ Filtering of weighted networks.
- ★ Entropic filtering of concepts.
- ★ Results with “*Special Effects*”.
 - Take home messages
 - Questions

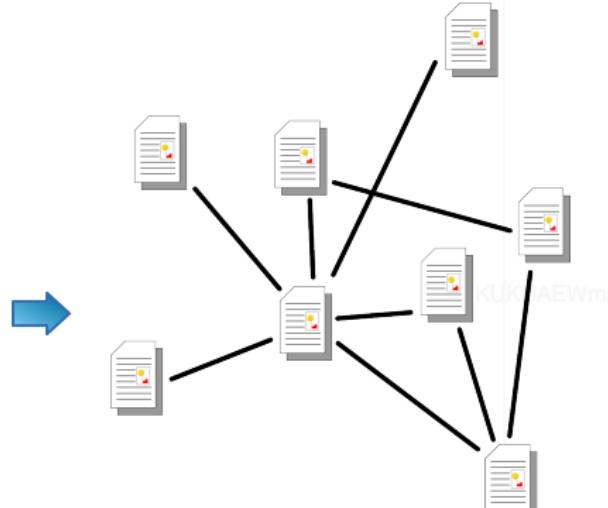
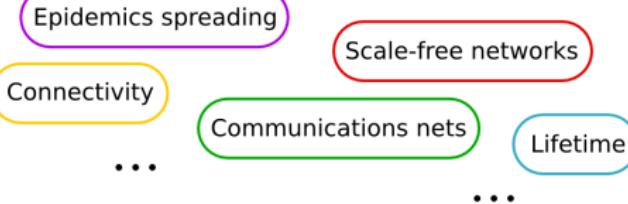
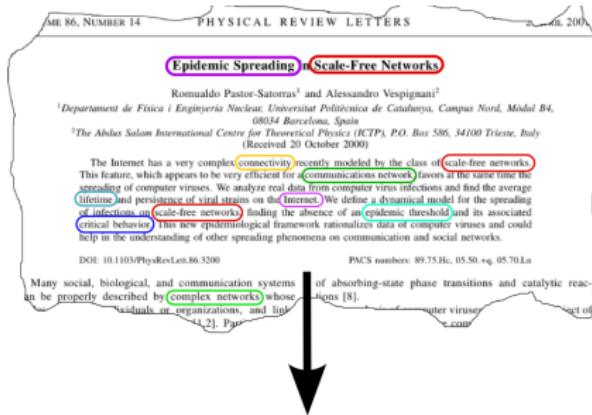
IAEWm

Section 1

Similarity networks

KUKUAEWm

Networks of similarity between papers



Similarity network

TF-IDF and similarity

α

1	1	1	0	1	0	0
---	---	---	---	---	---	---

$c_1 \ c_2 \ c_3 \ c_4 \ c_5 \ c_6 \ c_7$

$$TF-IDF_{\alpha c} = u_{\alpha c} = \underbrace{tf_{\alpha c}}_{TF} \underbrace{\log \left(\frac{N}{df_c} \right)}_{IDF}.$$

KUKUAEWm

TF-IDF and similarity

Edge weight

$$w_{\alpha\beta} = \frac{\vec{u}_\alpha \cdot \vec{u}_\beta}{\|\vec{u}_\alpha\| \|\vec{u}_\beta\|},$$

$$w_{\alpha\beta} \in [0, 1],$$

α	2	43	0	18	0	11	27
----------	---	----	---	----	---	----	----

β	13	5	9	0	30	0	0
---------	----	---	---	---	----	---	---

$C_1 \quad C_2 \quad C_3 \quad C_4 \quad C_5 \quad C_6 \quad C_7$

TF-IDF and similarity

Edge weight

$$w_{\alpha\beta} = \frac{\vec{u}_\alpha \cdot \vec{u}_\beta}{\|\vec{u}_\alpha\| \|\vec{u}_\beta\|},$$

$$w_{\alpha\beta} \in [0, 1],$$

$$\begin{aligned} w_{\alpha\beta} &= \frac{(13 \times 2) + (43 \times 5)}{55.02 \times 34.28} = \\ &= \frac{241}{1886.09} \simeq 0.13. \end{aligned}$$

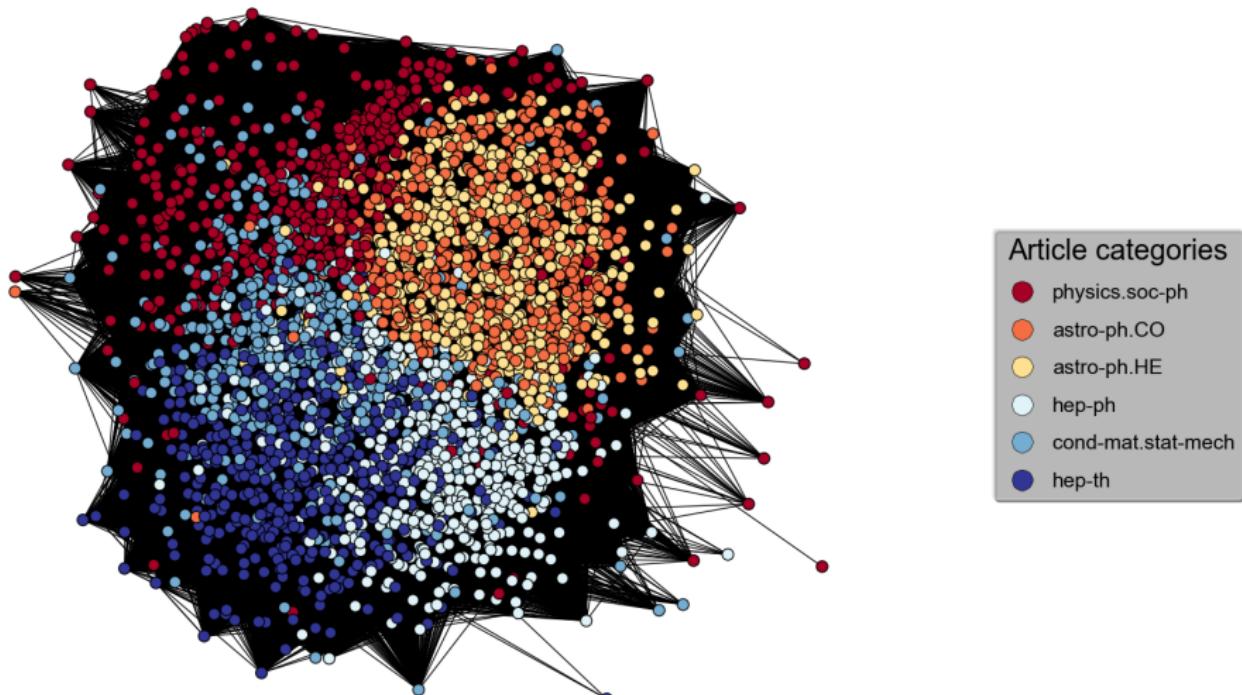
α	2	43	0	18	0	11	27
----------	---	----	---	----	---	----	----

β	13	5	9	0	30	0	0
---------	----	---	---	---	----	---	---

$C_1 \quad C_2 \quad C_3 \quad C_4 \quad C_5 \quad C_6 \quad C_7$

Network properties (Phys2013_pc)

N_c	N_a	$\langle k \rangle$	k_{max}	$\langle C \rangle$	d	$\langle L \rangle$	M	T
10661	52979	19333.522	46504	0.557	3	1.635	1	6





KEEP
CALM

AND

HOUSTON, WE

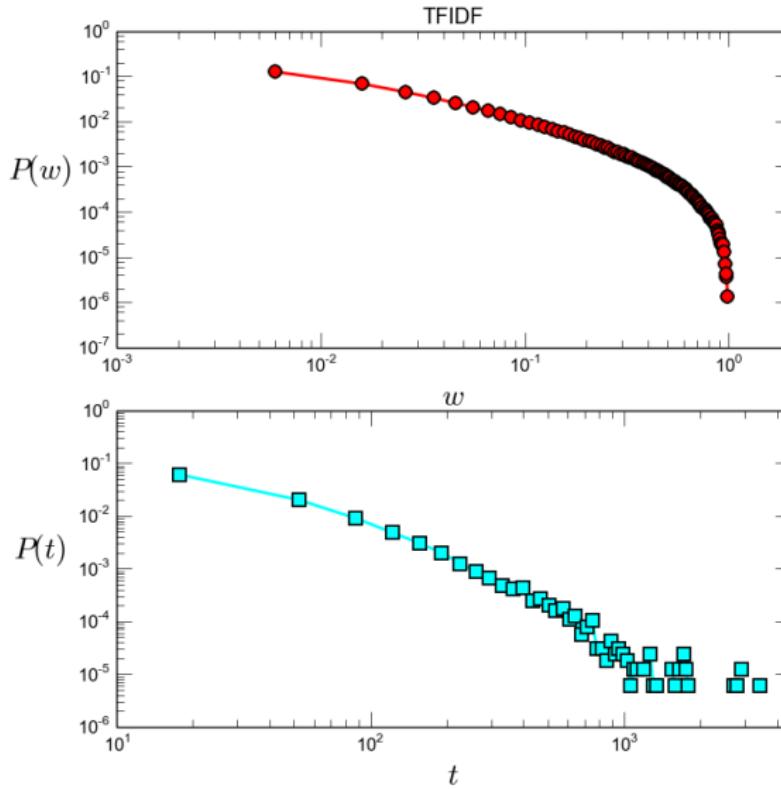
HAVE A PROBLEM

Section 2

Filtering

KUKUAEWm

Edge sparsification methods



KUKUAEWm

Edge sparsification methods

PHYSICAL REVIEW E

statistical, nonlinear, and soft matter physics

Highlights Recent Accepted Authors Referees Search About

Information filtering in complex weighted networks

Filippo Radicchi, José J. Ramasco, and Santo Fortunato

Phys. Rev. E **83**, 046101 – Published 1 April 2011

Article

References

Citing Articles (8)

PDF

HTML

Export Citation



ABSTRACT

Many systems in nature, society, and technology can be described as networks, where the vertices are the system's elements, and edges between vertices indicate the interactions between the corresponding elements. Edges may be weighted if the interaction strength is measurable. However, the full network information is often redundant because tools and techniques from network analysis

- Radicchi, F., et al. *Information filtering in complex weighted networks*. Physical Review E, **83** 046101. (2011).

Edge sparsification methods

Institution: EPFL
 Proceedings of the National Academy of Sciences of the United States of America

CURRENT ISSUE // ARCHIVE // NEWS & MULTIMEDIA // FOR AUTHORS // ABOUT PNAS // COLLECTED ARTICLES // BROWSE BY TOPIC // EARLY EDITION

[Home](#) > Current Issue > vol. 106 no. 16 > M. Ángeles Serrano, 6483–6488, doi: 10.1073/pnas.0808904106

 CrossMark
click for updates

Extracting the multiscale backbone of complex weighted networks

M. Ángeles Serrano^{a,1}, Marián Boguña^b and Alessandro Vespignani^{c,d}

Author Affiliations

Edited by Peter J. Bickel, University of California, Berkeley, CA, and approved March 2, 2009 (received for review September 9, 2008)

Abstract | Full Text | Authors & Info | Figures | SI | Metrics | Related Content |  |  +SI

This Issue

PNAS April 21, 2009
vol. 106 no. 16
Masthead (PDF)
Table of Contents

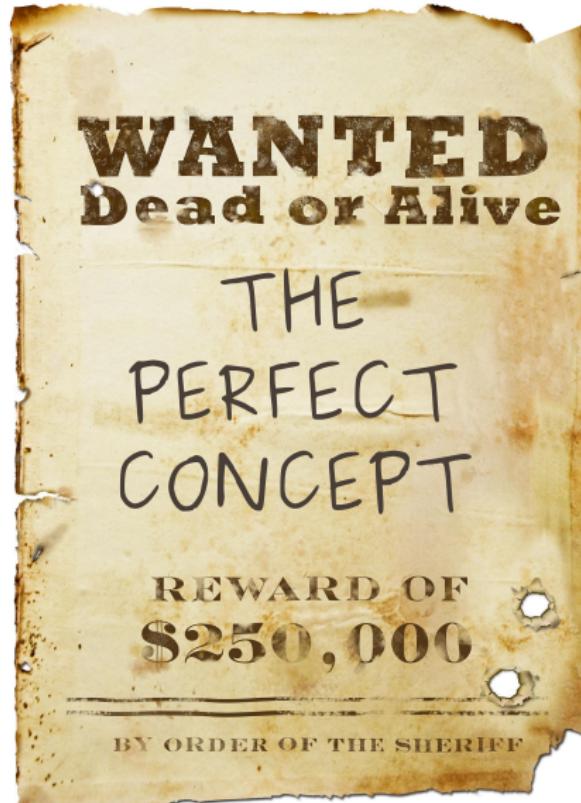
◀ PREV ARTICLE ▶ NEXT ARTICLE ▶

 View this article with LENS beta

Don't Miss

- Serrano M.A., et al. *Extracting the multiscale backbone of complex weighted networks*. Proc. Natl. Acad. Sci. (USA) 106 6483 (2009).

Relevant concepts



KUKUAEWm

Relevant concepts

Key features

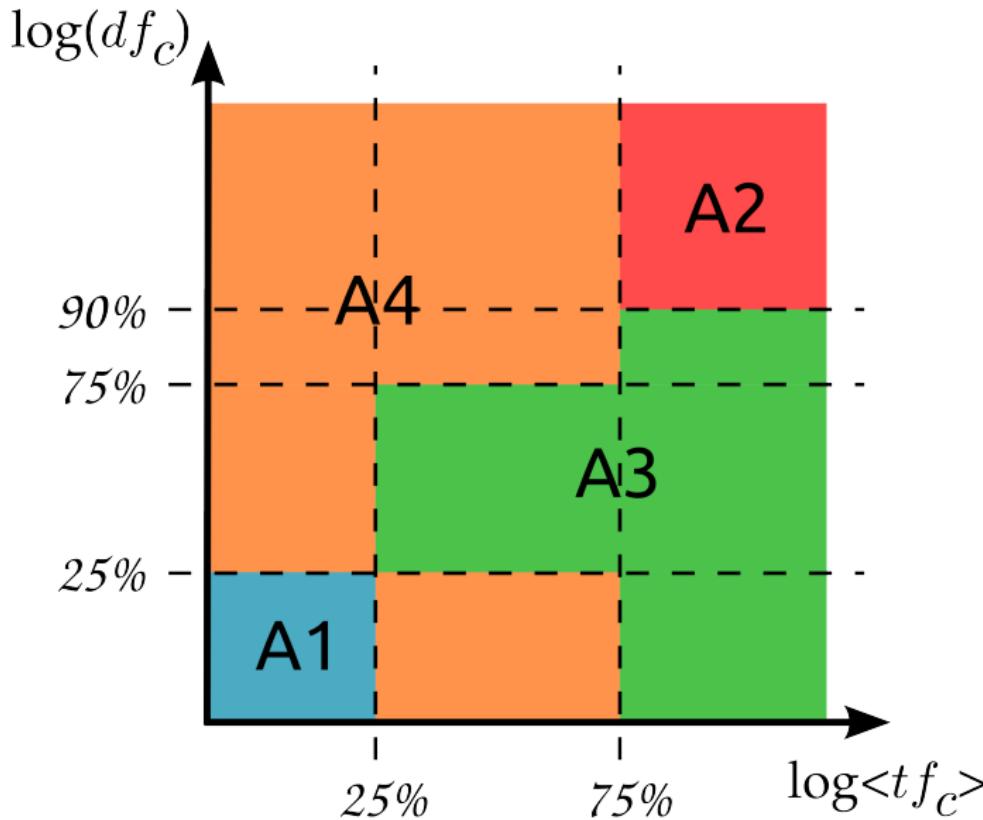
- # of papers a concept appears in

$df_c \rightarrow$ document frequency

- average # of times a concept appears inside a paper

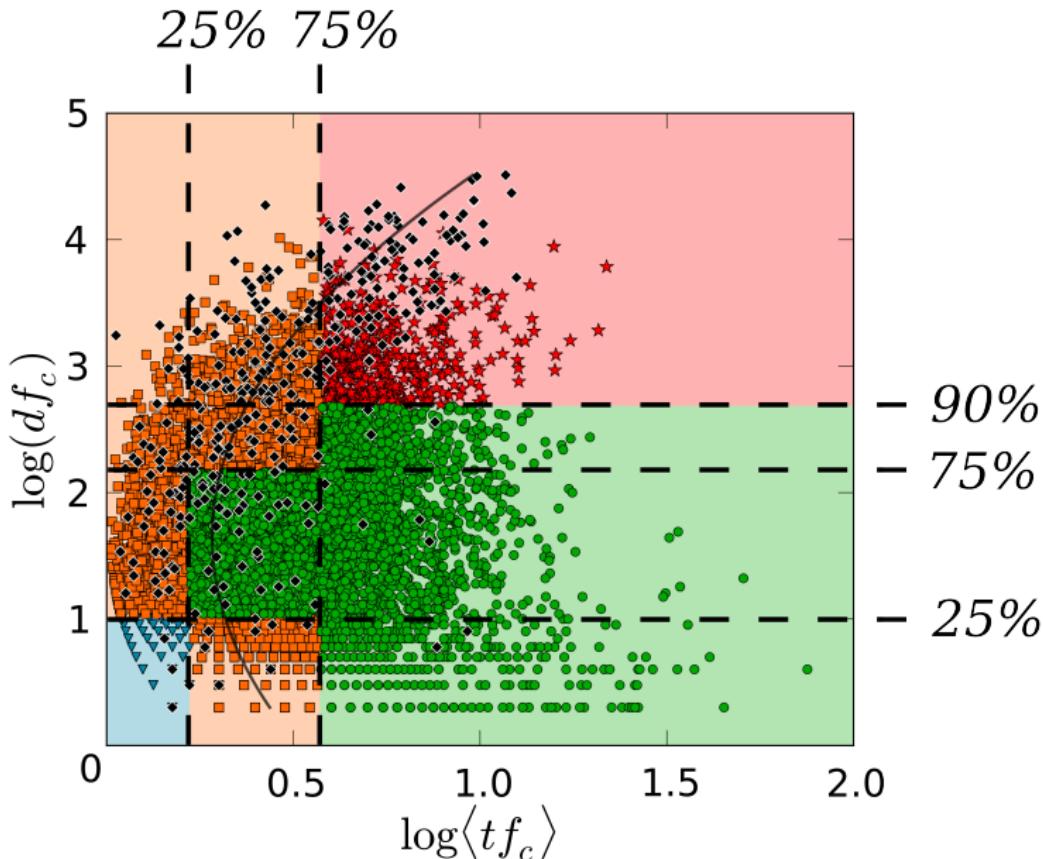
$\langle tf_c \rangle \rightarrow$ average term frequency

Bidimensional tessellation



KUKUAEWm

Bidimensional tessellation



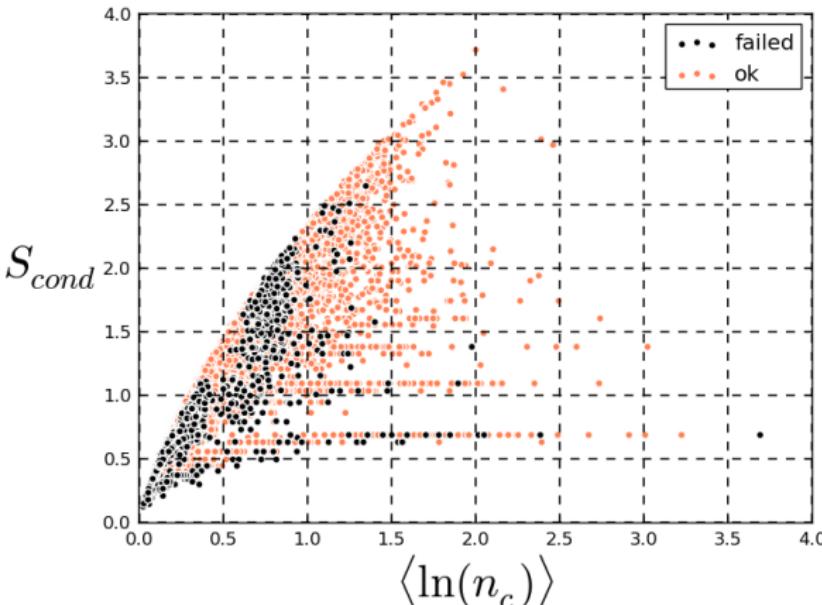
KUKUAEWm

Section 3

Entropic Filtering

KUKUAEWm

Maximum entropy



$$S = - \sum_{j=0}^{\infty} p_c(j) \ln p_c(j)$$

Maximum entropy

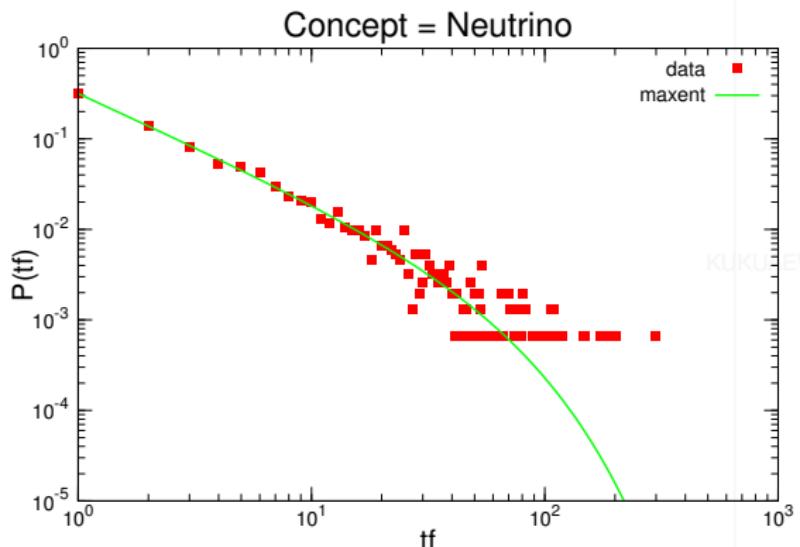
$$\sum_n p_n = 1$$

$$\sum_n p_n n = \langle n \rangle$$

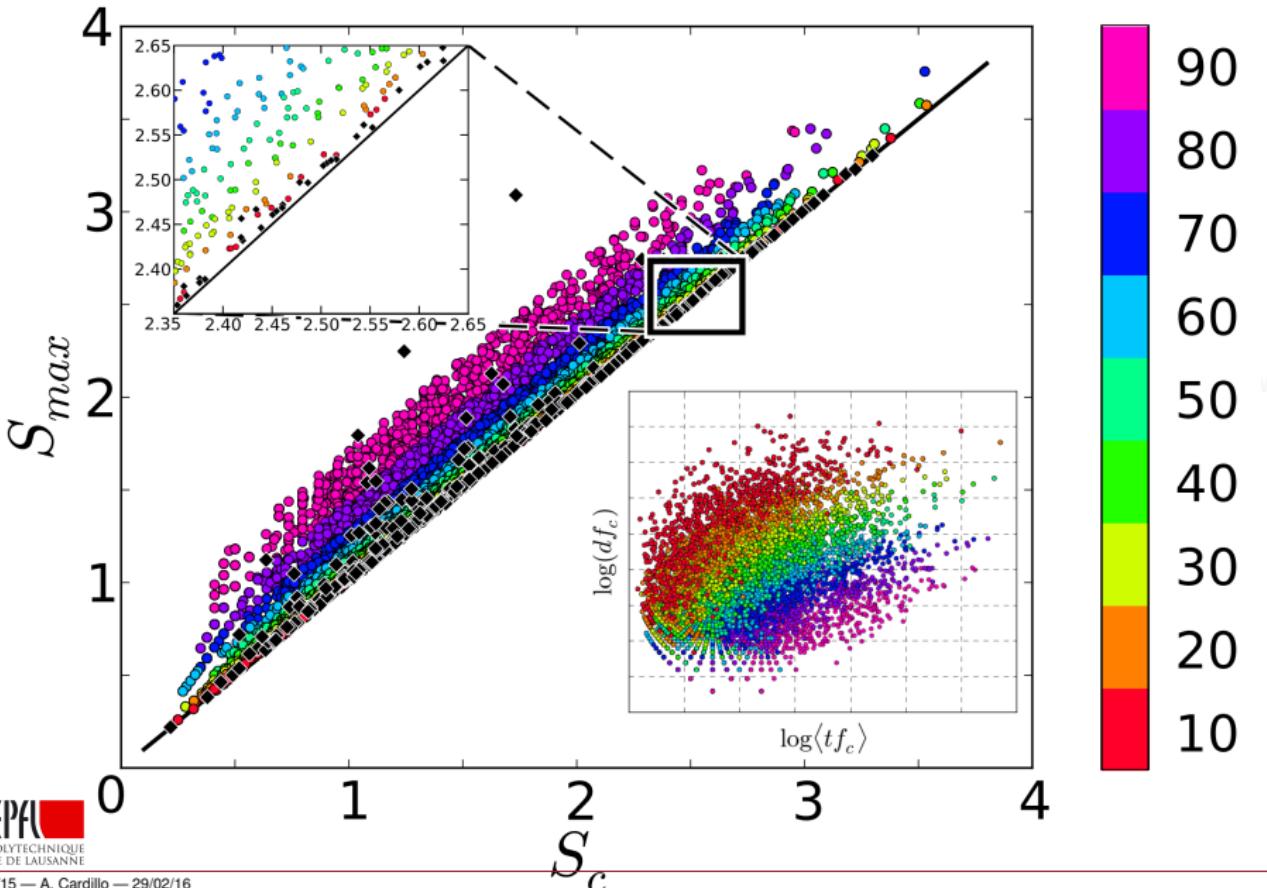
$$\sum_n p_n \ln n = \langle \ln n \rangle$$

$$\ln p_n + \lambda n + \mu \ln n = 0$$

$$p_n = \frac{1}{Z} e^{-\lambda n} n^{-\mu}$$



Maximum entropy



Section 4

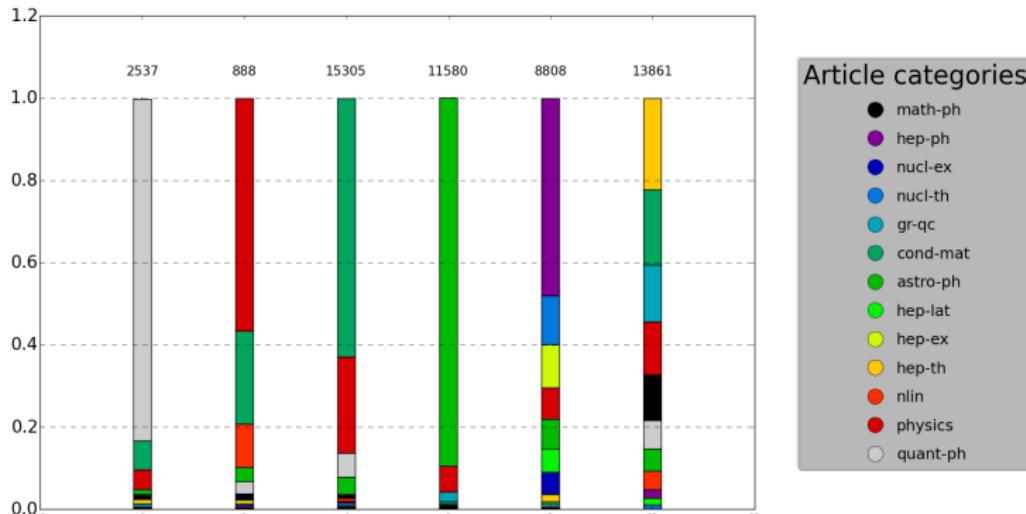
Results

KUKUAEWm

Topological properties

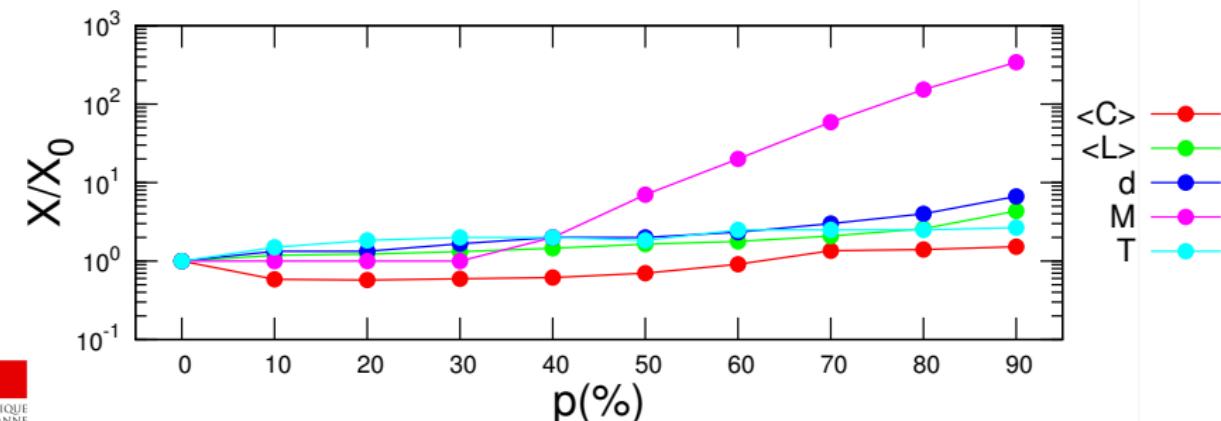
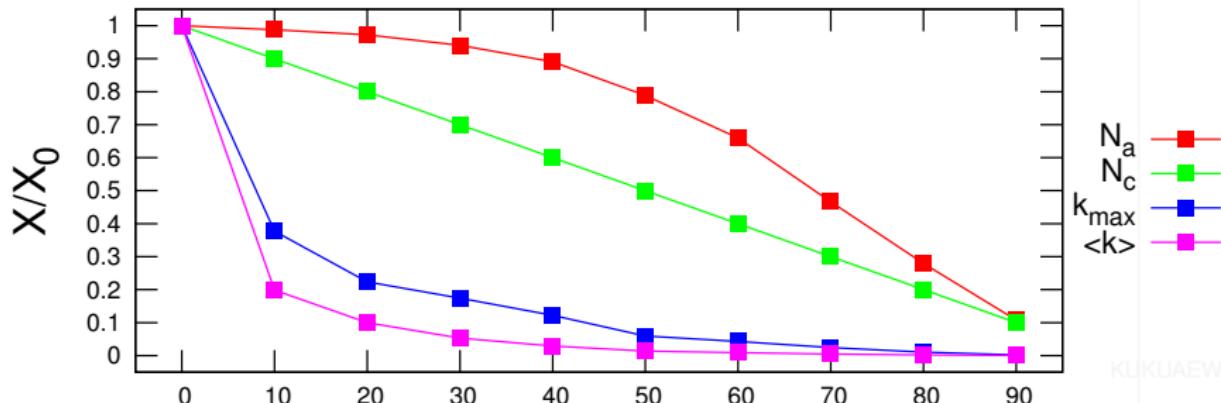
Original network

N_c	N_a	$\langle k \rangle$	k_{max}	$\langle C \rangle$	d	$\langle L \rangle$	M	T
10661	52979	19333.522	46504	0.557	3	1.635	1	6

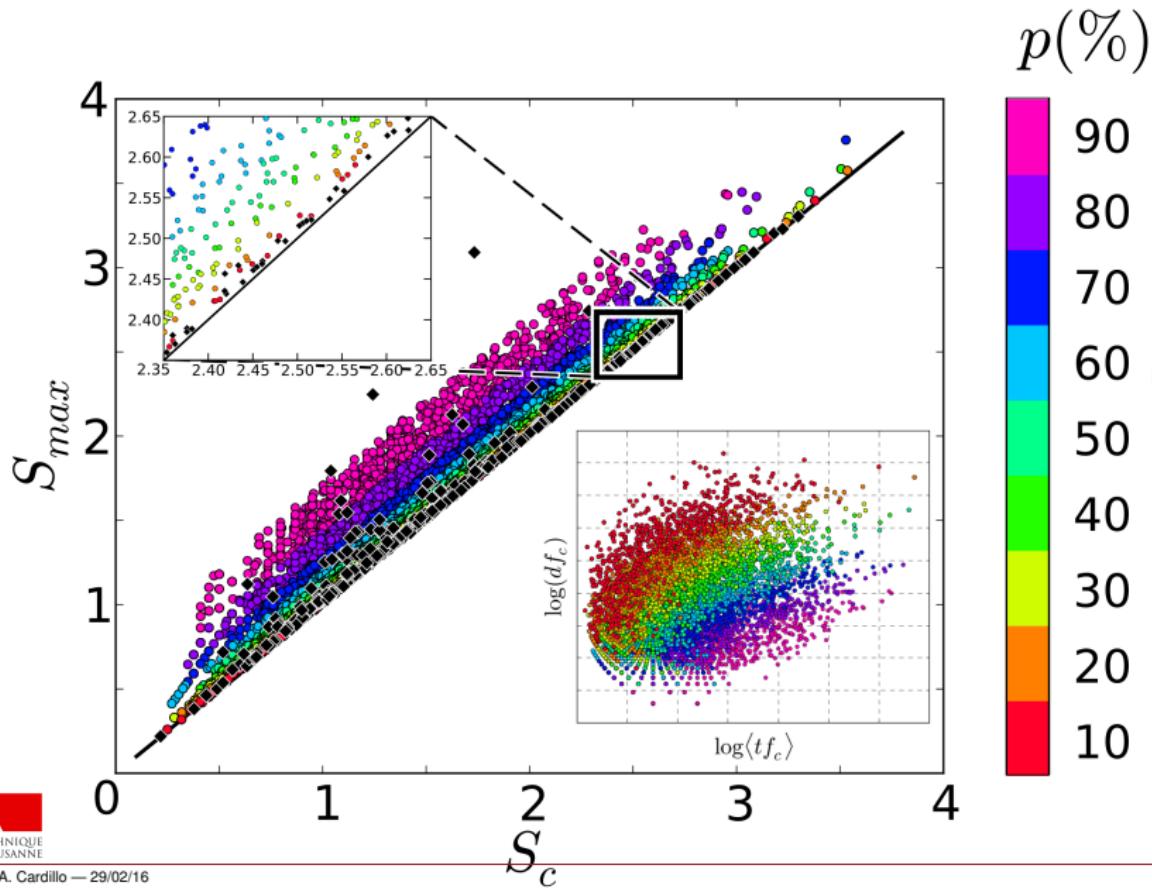


KUKUAEWm

Topological properties



What is a “generic concept”?



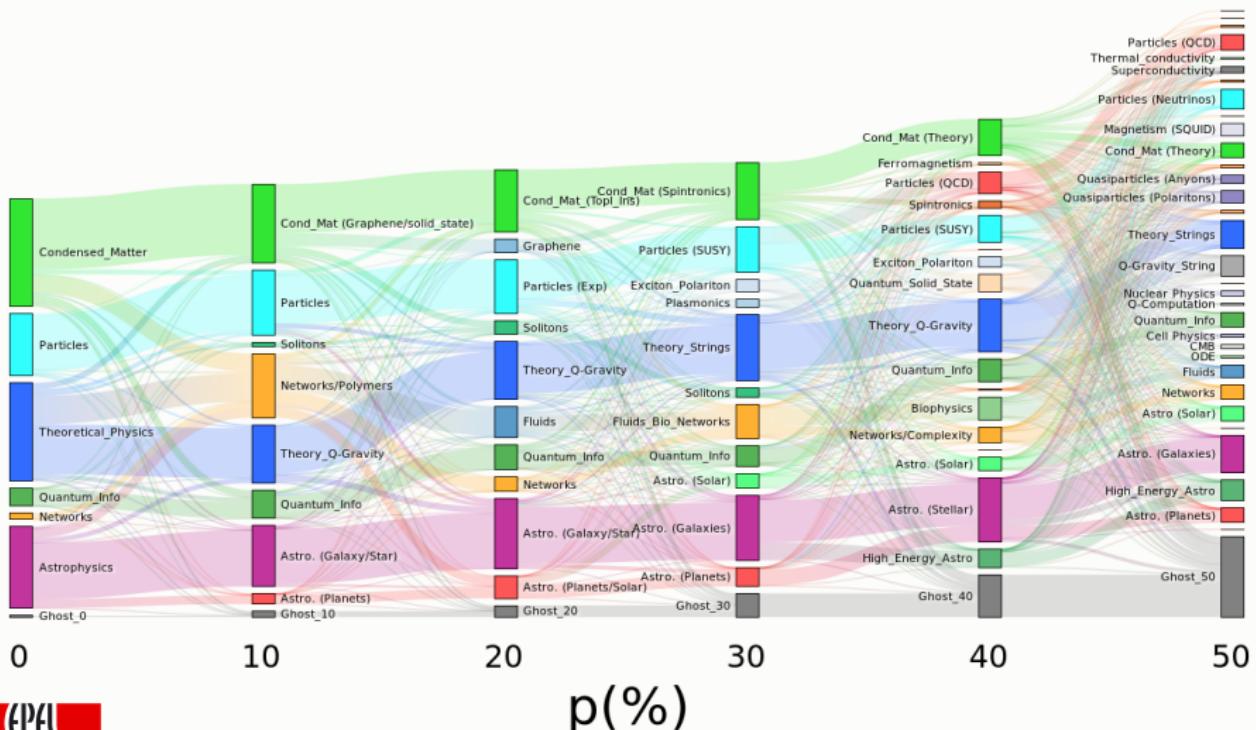
Community detection



▶ Link

Community detection

Data: Phys2013 $w_{min} = 0.01$



Section 5

Conclusions

KUKUAEWm

Summing up . . .

Take home messages

- We have used the maximum entropy principle to build a method to filter networks of similarity between documents.

IAEWm

Summing up . . .

Take home messages

- We have used the maximum entropy principle to build a method to filter networks of similarity between documents.
- The removal of common concepts allows to retrieve a more well defined structure of documents into topics.

Summing up . . .

Take home messages

- We have used the maximum entropy principle to build a method to filter networks of similarity between documents.
- The removal of common concepts allows to retrieve a more well defined structure of documents into topics.
- The method allows to identify collection dependent “*relevant concepts*” without requiring user validation.

Summing up . . .

What's next? Open questions

- Apply the methodology recursively.

JAEWm

Summing up . . .

What's next? Open questions

- Apply the methodology recursively.
- Validate the method applying it on collection of documents different from scientific papers.

IAEWm

Summing up . . .

What's next? Open questions

- Apply the methodology recursively.
- Validate the method applying it on collection of documents different from scientific papers.
- Study the evolution in time of “generality”.

Summing up ...

Reference

A. Martini *et al.* , *Automatic selection of relevant concepts in scientific publications* – in preparation

<http://bifi.es/~cardillo/>

alessio.cardillo@epfl.ch