



# 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/>



Once upon a time . . .

# Once upon a time ...



KUKUAEWm

# Once upon a time . . .



Nowdays ...



# Nowdays ...



KUAEWm

# Flood of information

## 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

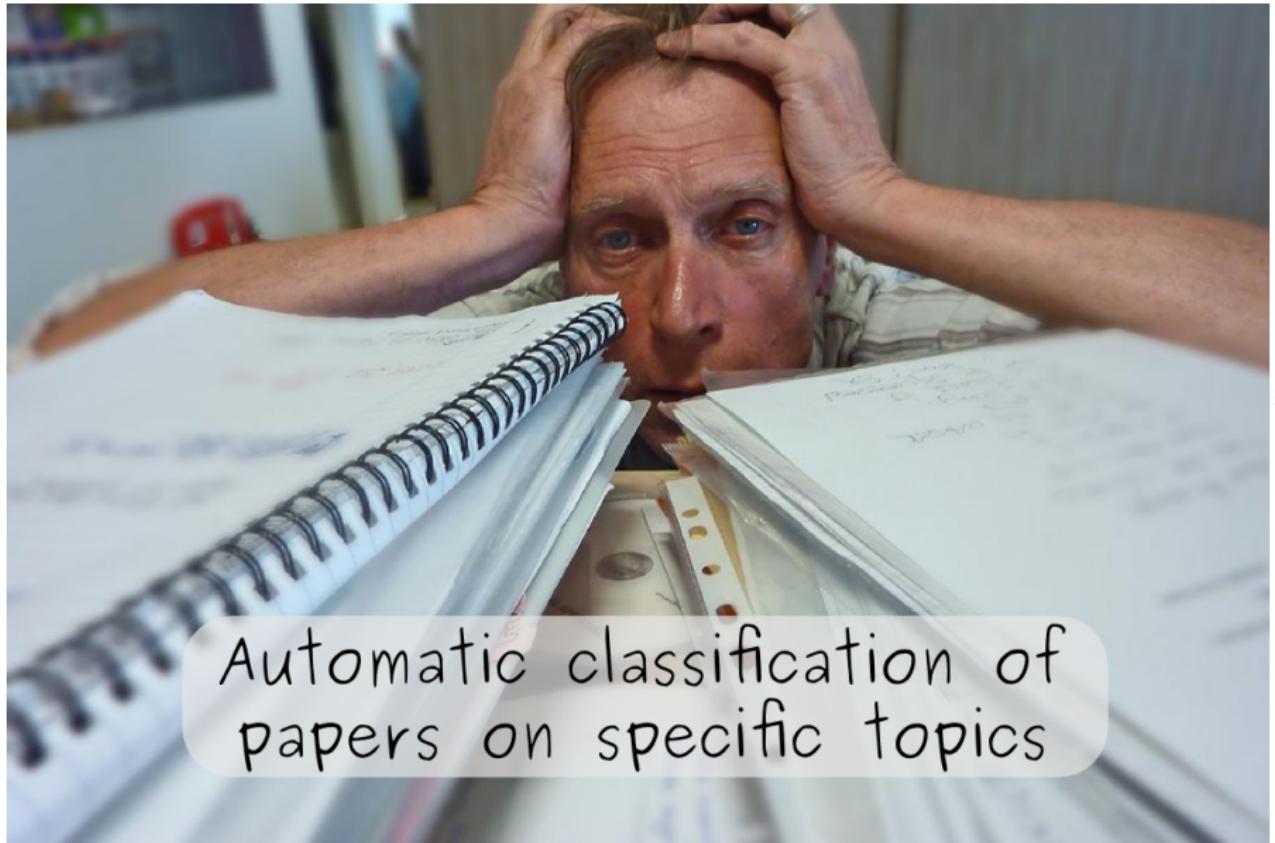
**German research agencies condemn animal-rights attack on neuroscientist**

### NEWS BLOG

## Global scientific output doubles every nine years

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

# Flood of information



Automatic classification of  
papers on specific topics

# Flood of information



International weekly journal of science

[Home](#) | [News & Comment](#) | [Research](#) | [Careers & Jobs](#) | [Current Issue](#) | [Archive](#) | [Audio & Video](#) | [For](#)

[Archive](#) > [Volume 513](#) > [Issue 7516](#) > [Toolbox](#) > [Article](#)

NATURE | TOOLBOX



## How to tame the flood of literature

Recommendation services claim to help researchers keep up with the most important papers without becoming overwhelmed.

Elizabeth Gibney

03 September 2014

# What do we need?

*There is an inherent problem to giving you information that you weren't actively searching for. **It has to be relevant** – so that we are not wasting your time – **but not too relevant**, because you already know about those articles.*

Anurag Acharya  
Google Scholar co-creator

## What do we need?

*There is an inherent problem to giving you information that you weren't actively searching for. **It has to be relevant** – so that we are not wasting your time – **but not too relevant**, because you already know about those articles.*

Anurag Acharya  
Google Scholar co-creator

*Semantic Scholar offers a few innovative features, including picking out the **most important keywords and phrases** from the text without relying on an author or publisher to key them in. “**It’s surprisingly difficult for a system to do this,**”*

Oren Etzioni  
CEO of AI2 (Semantic Scholar)

# What do we need?

**ScienceWISE**   [Ontology](#)   [Bookmarks](#)   [New articles](#)   [News](#)   [Introduction](#)   [Logout](#)

[Physics](#)   [Life Sciences beta](#)   [Digital Humanities](#)   [Information Technologies](#)

**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 physics  
 ...

**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

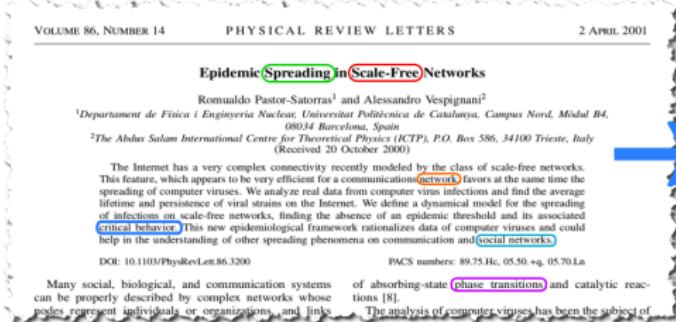
IAEWm

# Section 1

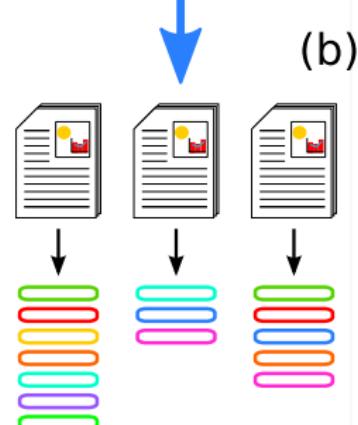
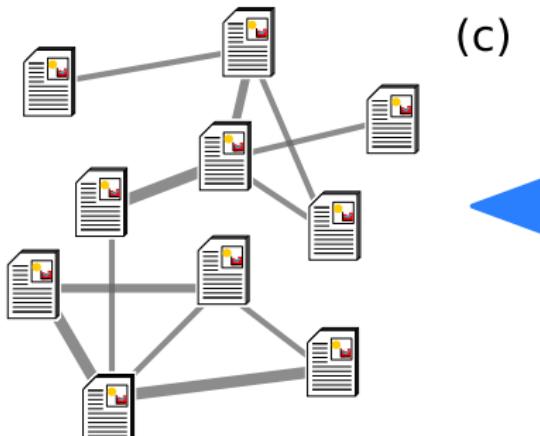
## Similarity networks

KUKUAEWm

# Networks of documents



(a) Scale-Free  
Social network  
Phase transition  
Network  
Spreading



# TF-IDF and similarity

$\alpha$	X	X	X	0	X	0	0
	$C_1$	$C_2$	$C_3$	$C_4$	$C_5$	$C_6$	$C_7$

KUKUAEWm

$$TF-IDF_{\alpha c} = u_{\alpha c} = \underbrace{tf_{\alpha c}}_{TF} \underbrace{\log \left( \frac{1}{df_c} \right)}_{IDF} = tf_{\alpha c} \log \left( \frac{N}{N_c} \right).$$

# TF-IDF and similarity

## Edge weight

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

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

$\beta$	2	43	0	18	0	11	27
$\alpha$	13	5	9	0	30	0	0
	$C_1$	$C_2$	$C_3$	$C_4$	$C_5$	$C_6$	$C_7$

# TF-IDF and similarity

## Edge weight

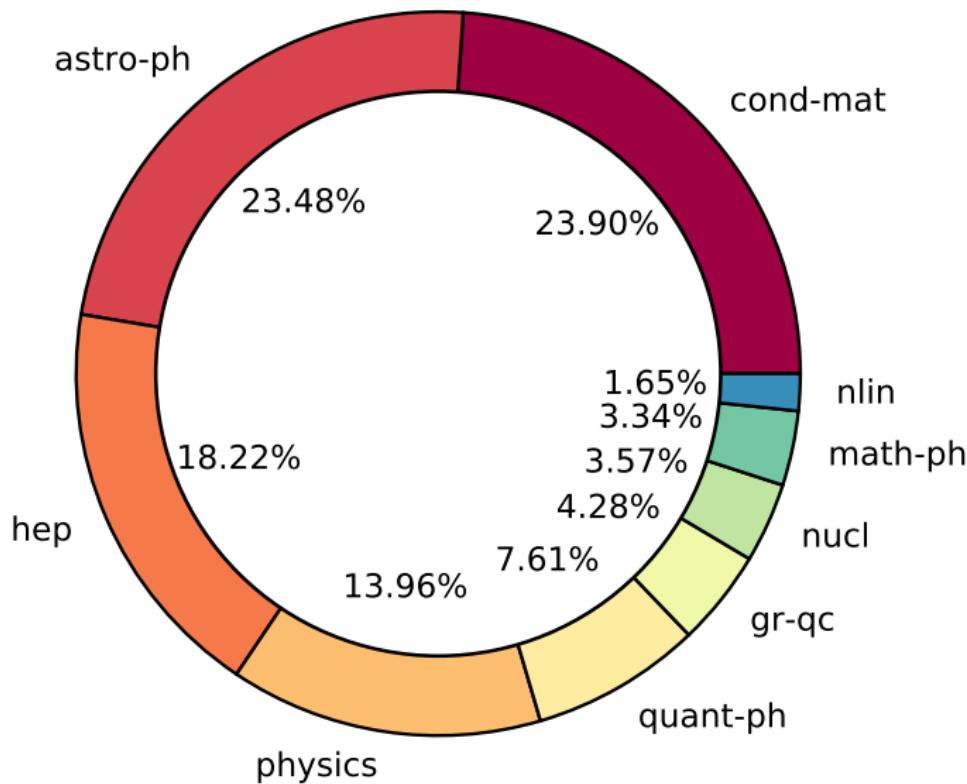
$$w_{\alpha\beta} = \frac{\vec{u}_\alpha \cdot \vec{u}_\beta}{\|\vec{u}_\alpha\| \|\vec{u}_\beta\|},$$

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

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

$\beta$	2	43	0	18	0	11	27
$\alpha$	13	5	9	0	30	0	0
	$C_1$	$C_2$	$C_3$	$C_4$	$C_5$	$C_6$	$C_7$

# The data: 2013 Physics arXiv

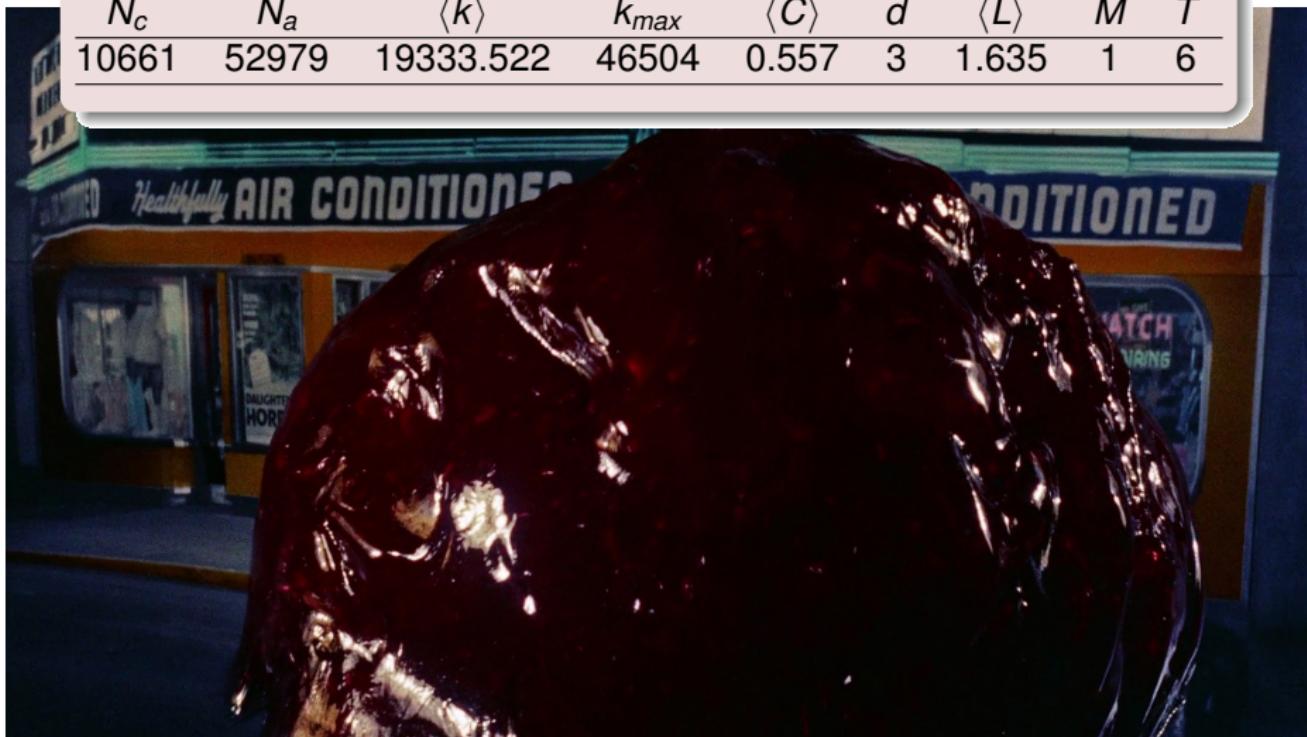


KUAEWm

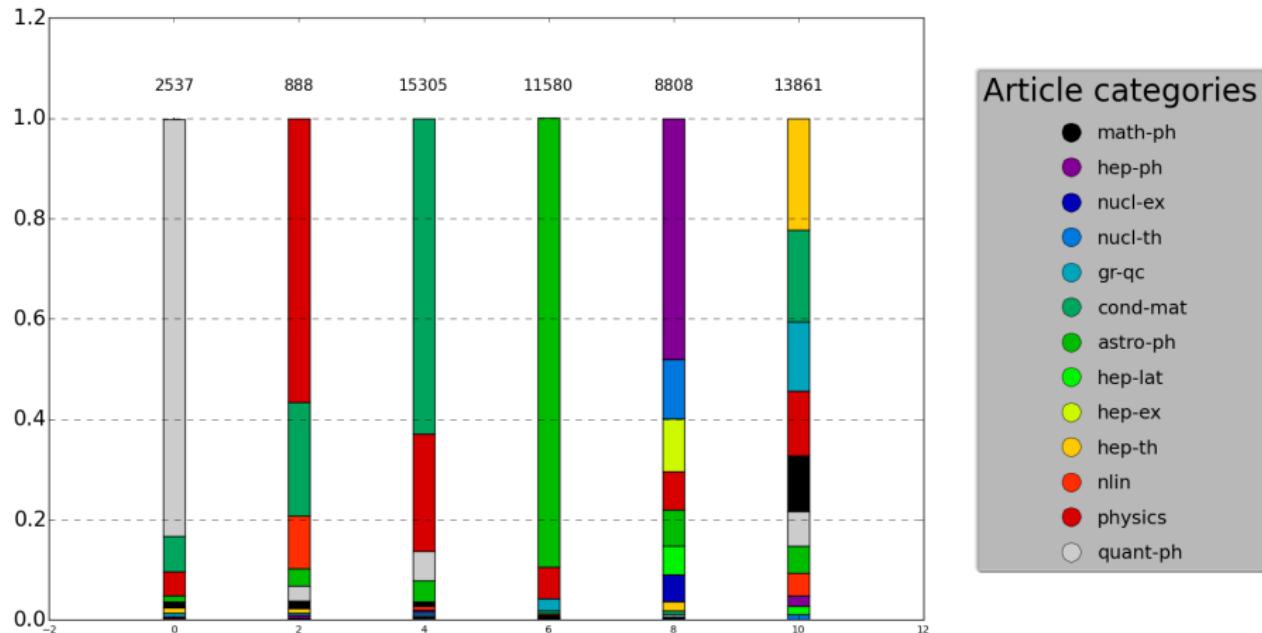
# The data: 2013 Physics arXiv

## Network properties

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



# The data: 2013 Physics arXiv



$$\rho = \frac{K}{K_{\max}} \simeq 36\%$$

# The data: 2013 Physics arXiv



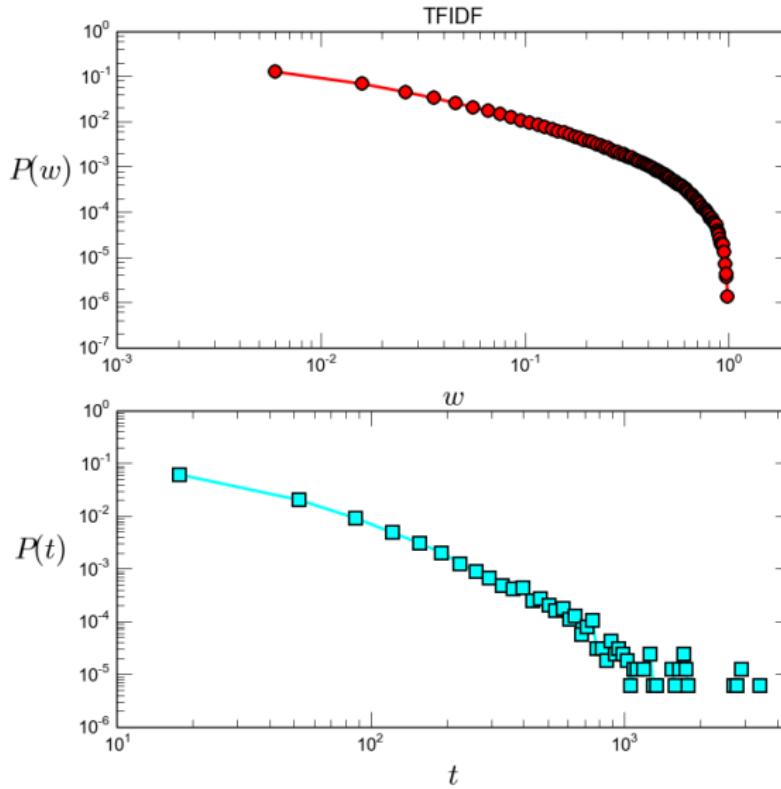
KUKUAEWm

## Section 2

### Filtering

KUKUAEWm

# Edge pruning/sparsification methods



KUKUAEWm

# Edge pruning/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<sup>a,1</sup>, Marián Boguña<sup>b</sup> and Alessandro Vespignani<sup>c,d</sup>

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).

# Edge pruning/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

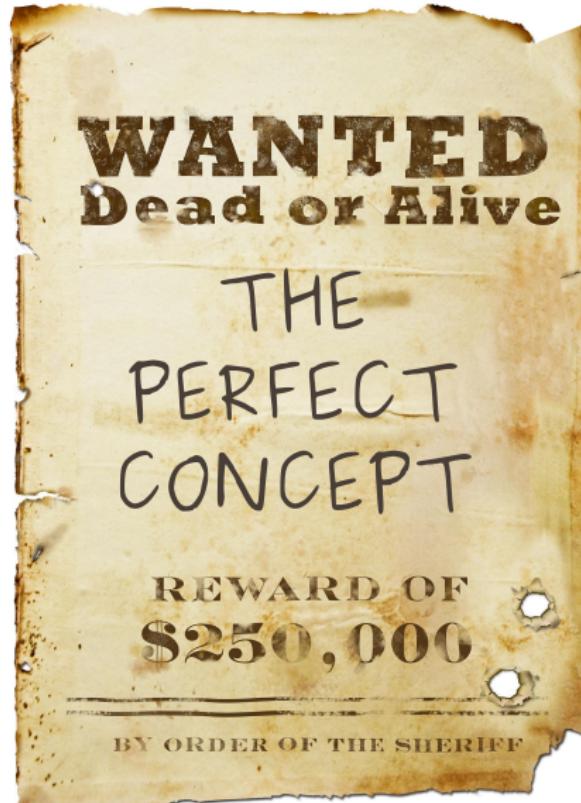


#### ABSTRACT

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).

# Relevant concepts



KUKUAEWm

# Relevant concepts

## Key features

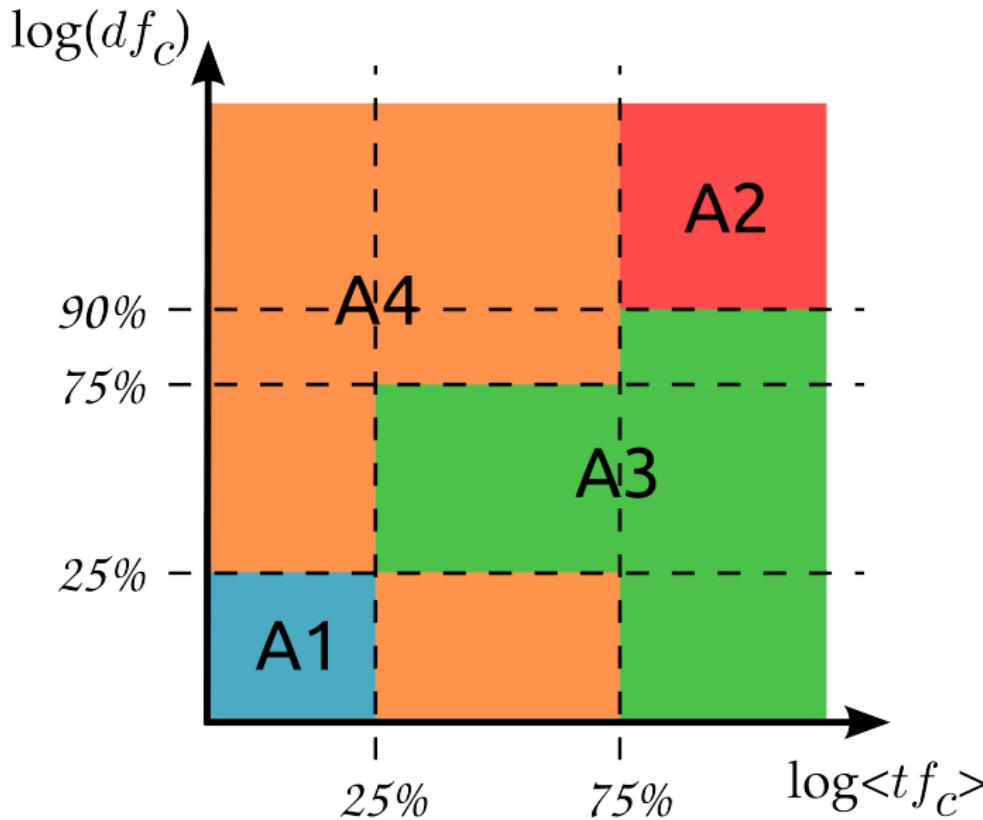
- # of papers a concept appears in

$df_c \rightarrow$  document frequency

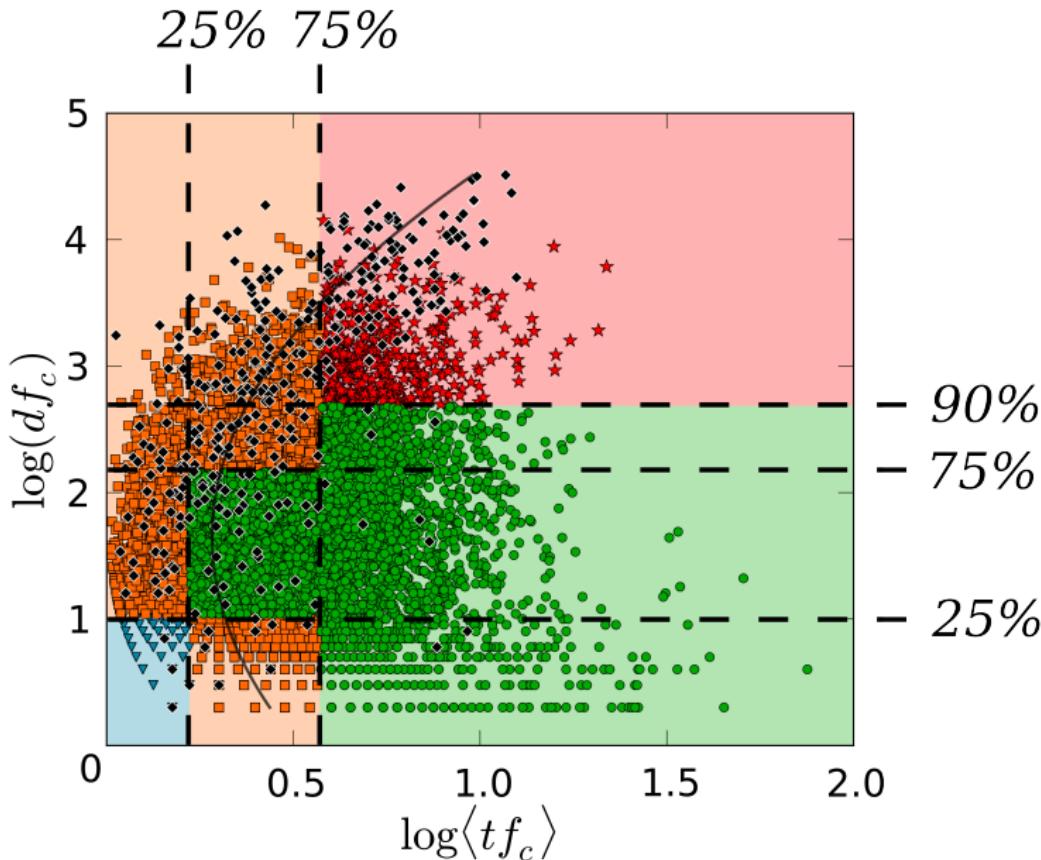
- average # of times a concept appears inside a paper

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

# Bidimensional tessellation



# Bidimensional tessellation



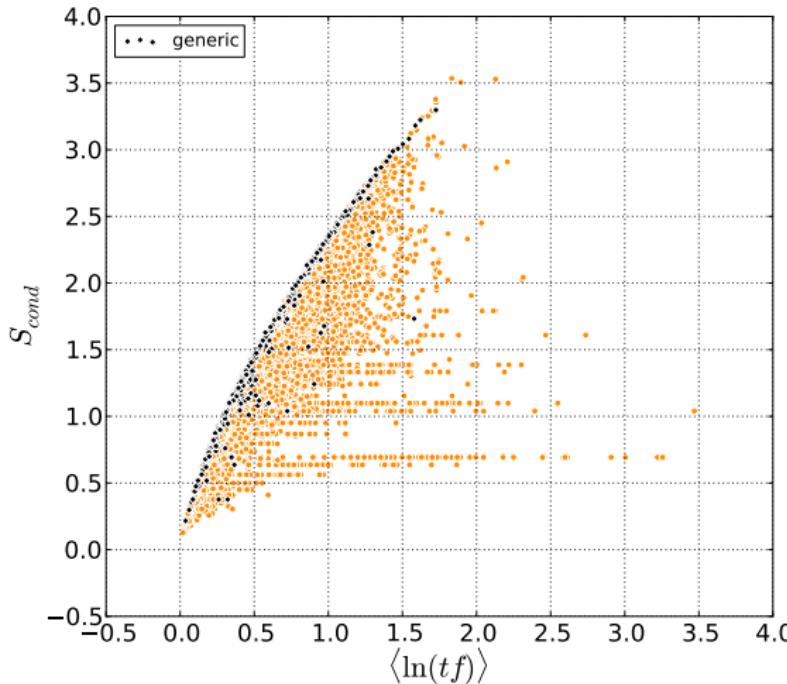
KUKUAEWm

## Section 3

### Entropic Filtering

KUKUAEWm

# Maximum entropy



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

# Maximum entropy

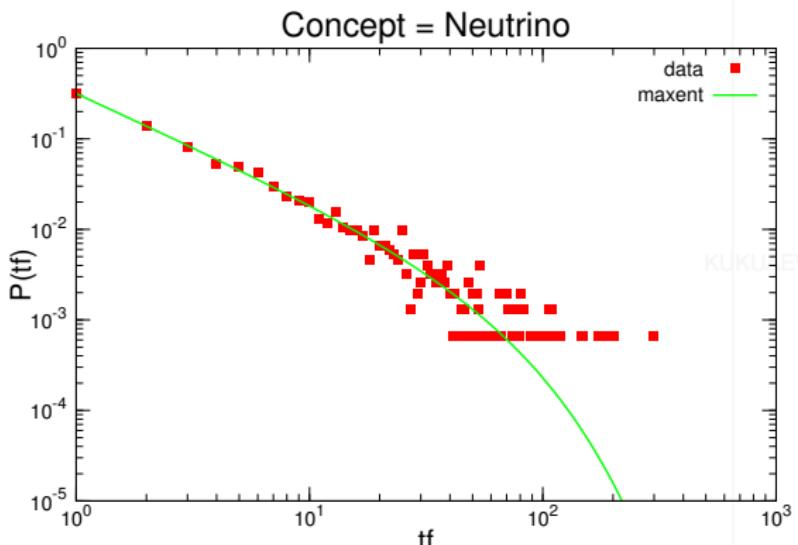
$$\sum_n p_n = 1$$

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

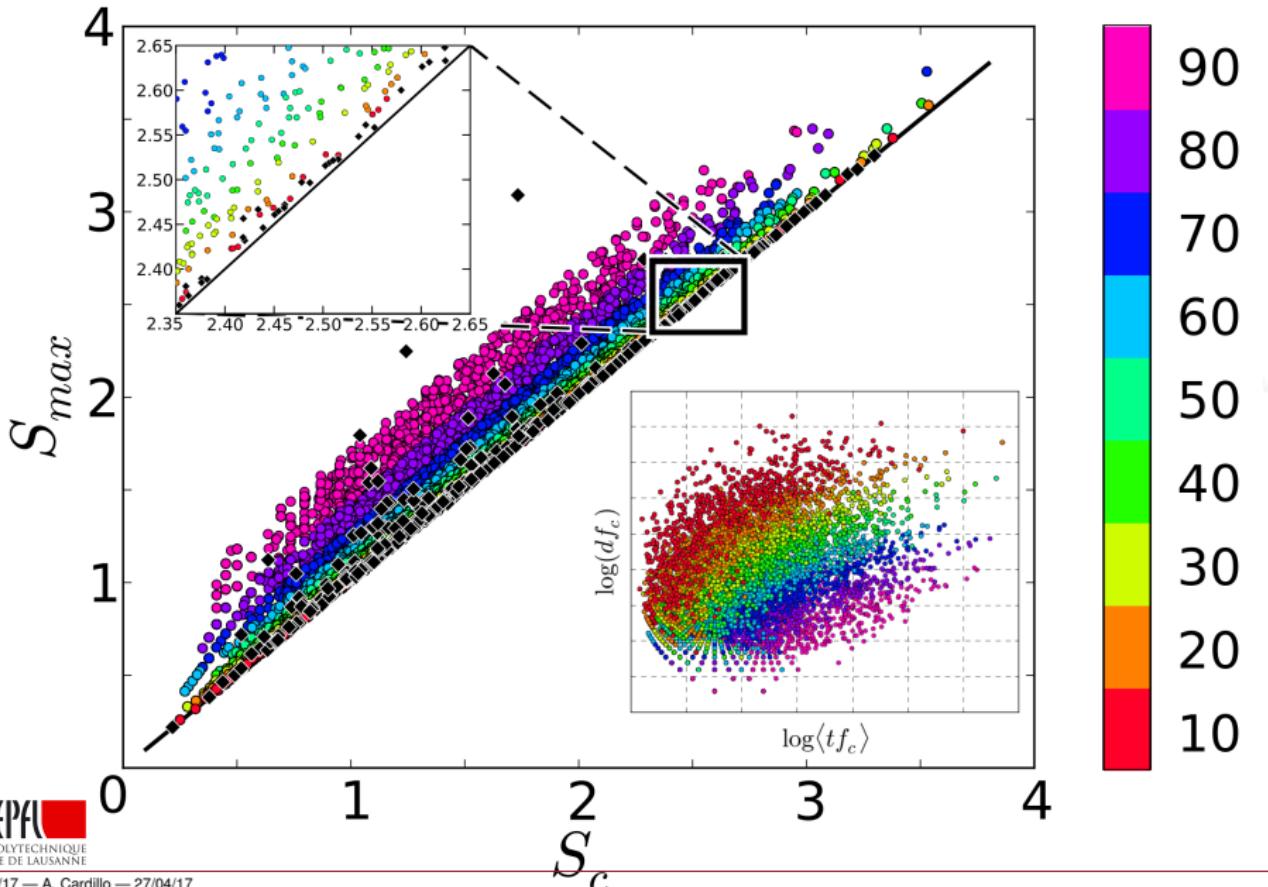
$$\sum_n p_n \ln n = \langle \ln n \rangle$$

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

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



# Maximum entropy



## Section 4

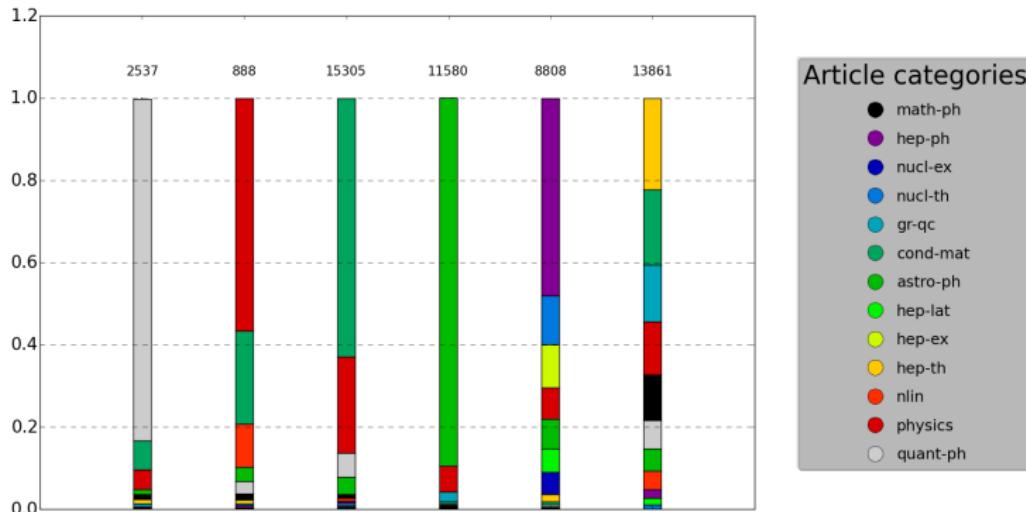
### Results

KUKUAEWm

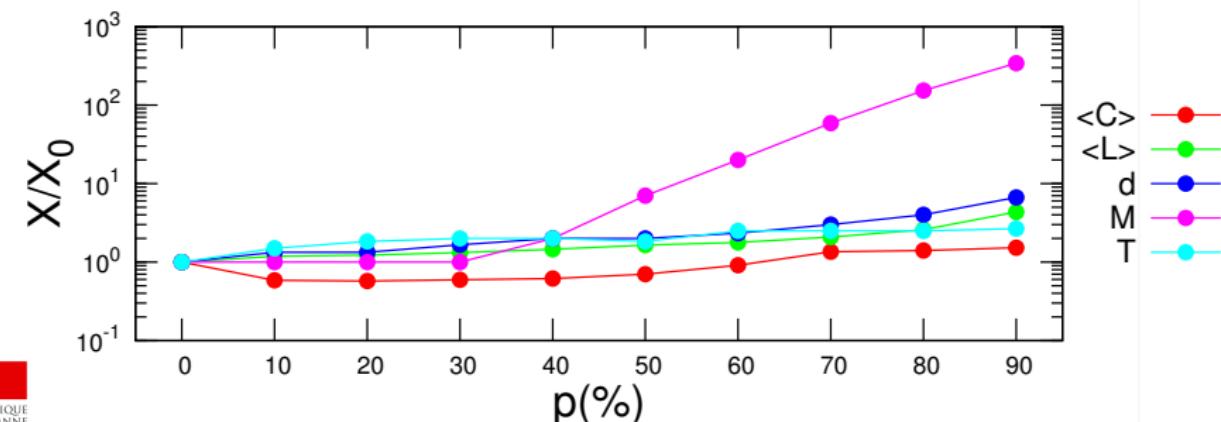
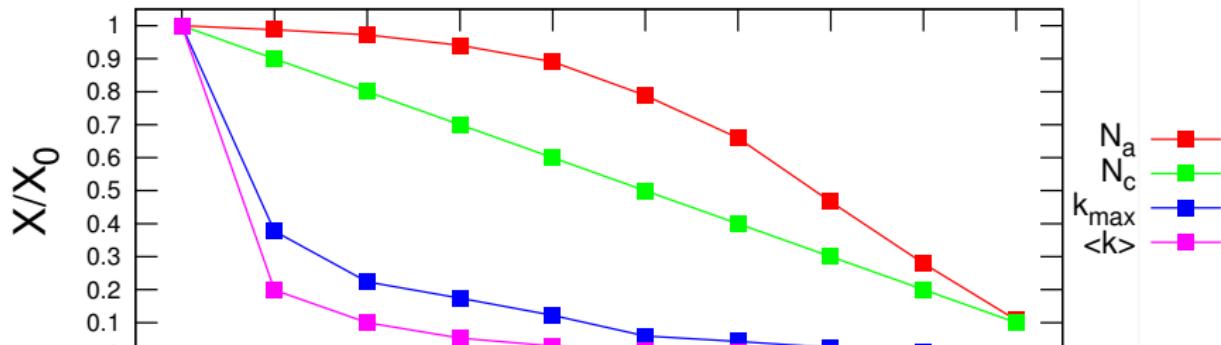
# Topological properties

## Original network

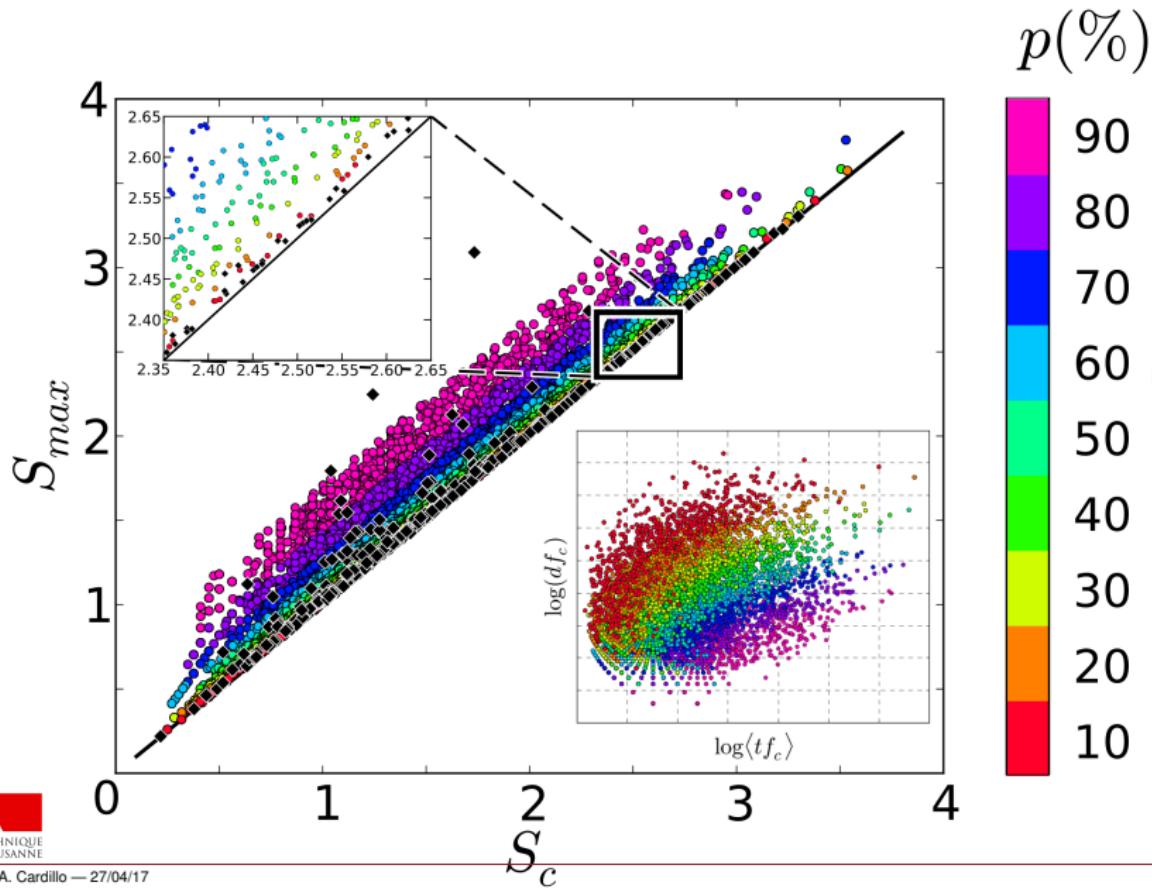
$N_c$	$N_a$	$\langle k \rangle$	$k_{max}$	$\langle C \rangle$	$d$	$\langle L \rangle$	$M$	$T$
10661	52979	19333.522	46504	0.557	3	1.635	1	6



# Topological properties



# What is a “generic concept”?



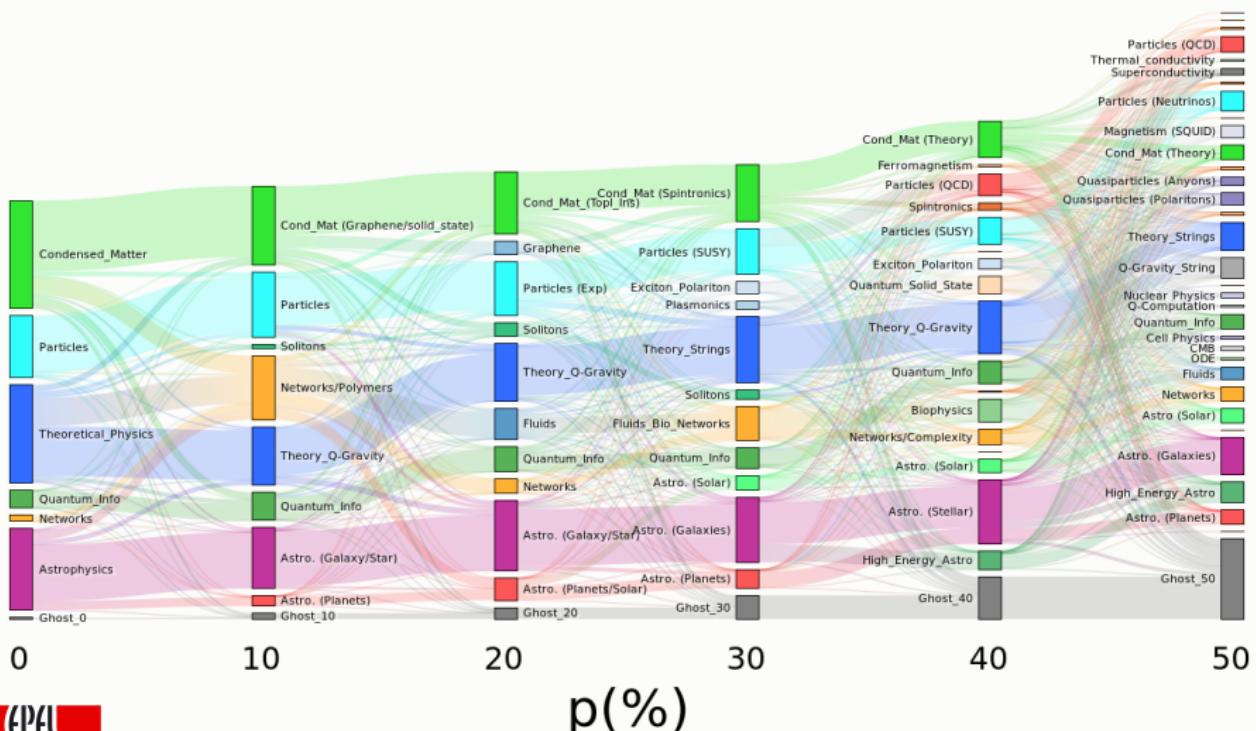
# Community detection



▶ 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 method allows to identify collection dependent “*relevant concepts*” without requiring user validation.

# Summing up . . .

## Take home messages

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

# Summing up . . .

## What's next? Open questions

- Study the evolution in time of “generality”.

JAEWm

# Summing up . . .

## What's next? Open questions

- Study the evolution in time of “generality”.
- Study the relation between concepts.

IAEWm

# Summing up . . .

## What's next? Open questions

- Study the evolution in time of “generality”.
- Study the relation between concepts.
- Use the method to build ontologies.

IAEWm

# Summing up ...

## Reference

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

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

alessio.cardillo@epfl.ch

# Acknowledgements



***Paolo De Los Rios***



***Andrea Martini***



***Alex Constantin***  
**&**  
***ScienceWISE team***

# Acknowledgements



FONDS NATIONAL SUISSE  
SCHWEIZERISCHER NATIONALFONDS  
FONDO NAZIONALE SVIZZERO  
SWISS NATIONAL SCIENCE FOUNDATION

KUKUAEWm

Academic freedom matters.  
**#IStandWithCEU**

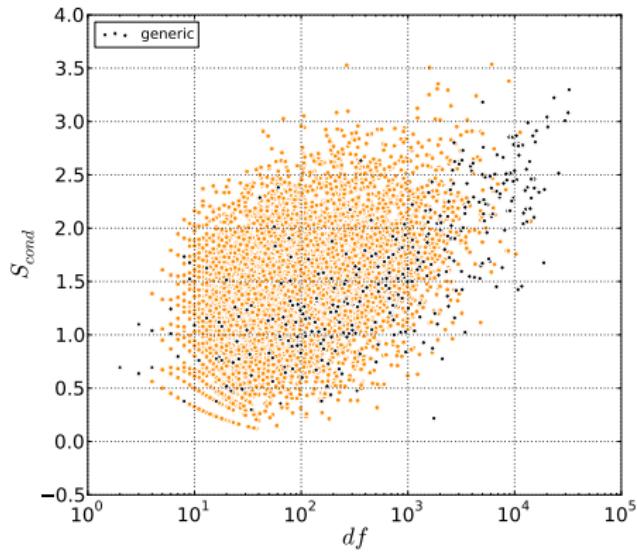
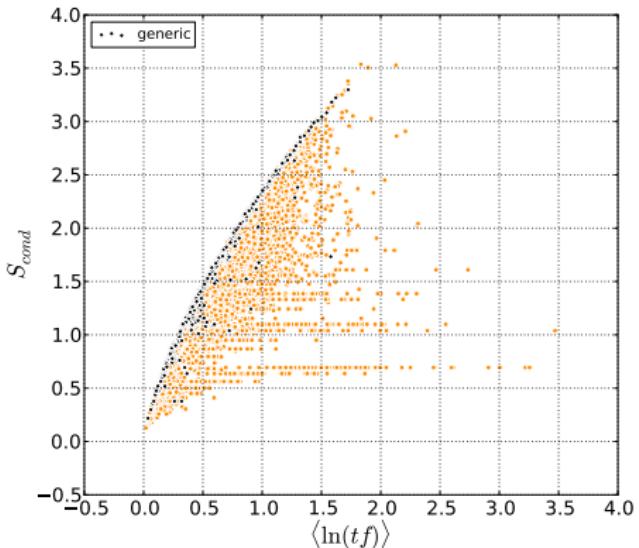
KUAEWm

# Part II

## Appendix

KUKUAEWm

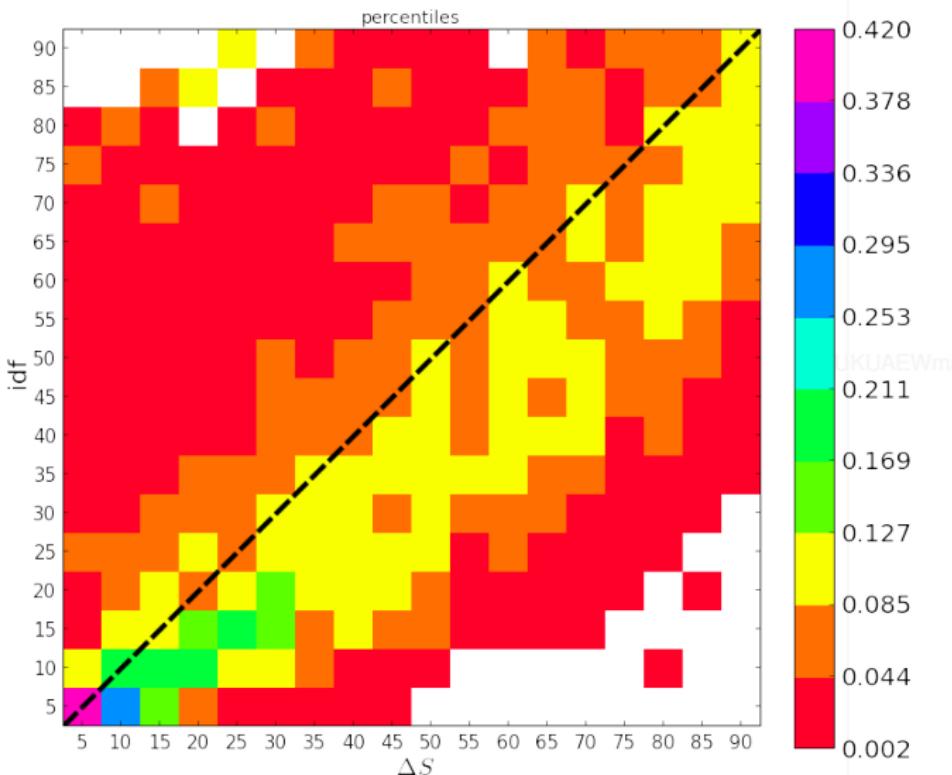
# Why not $df$ ?



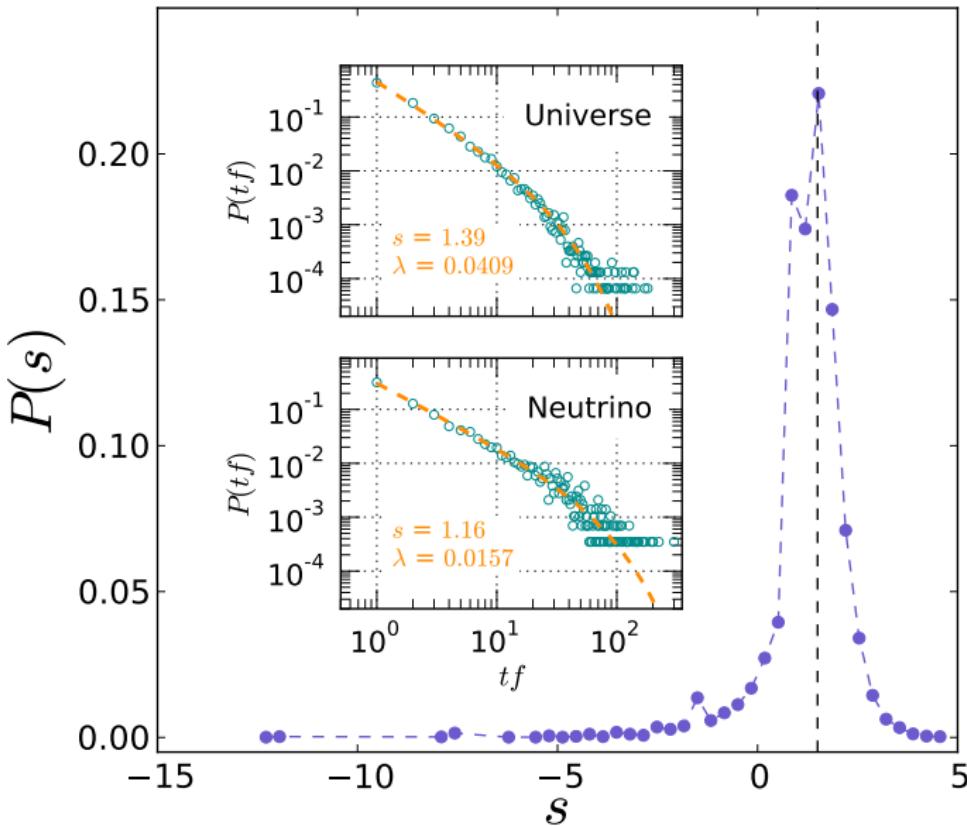
# Why not $df$ ?

Jaccard Score

$$J = \frac{|A \cap B|}{|A \cup B|}.$$

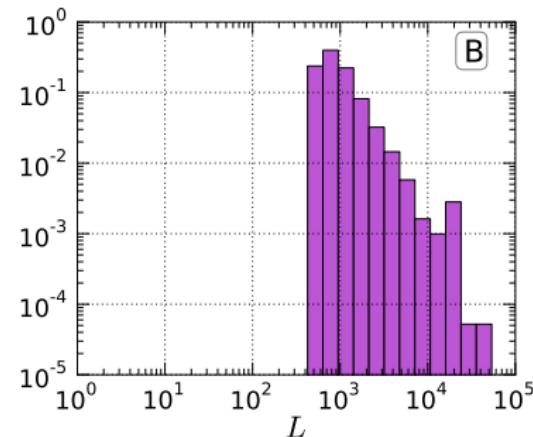
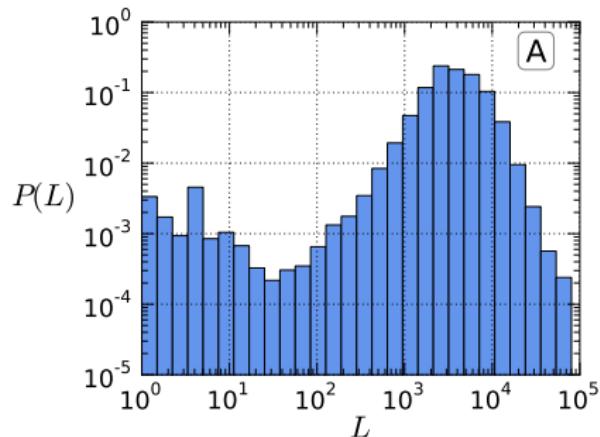


# Maximum entropy – why power law?



KUKUAEWm

# Maximum entropy – TF density

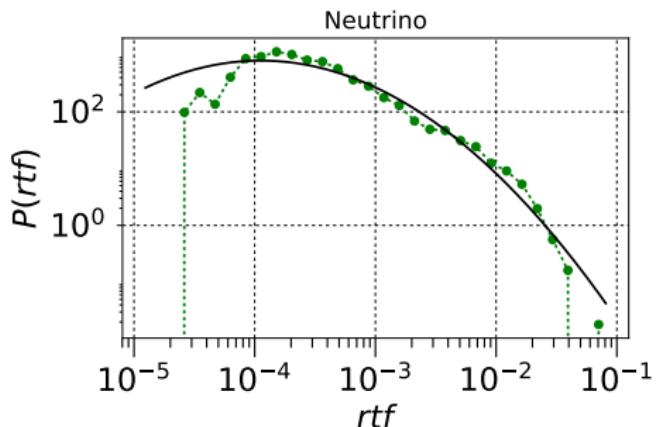


# Maximum entropy – TF density

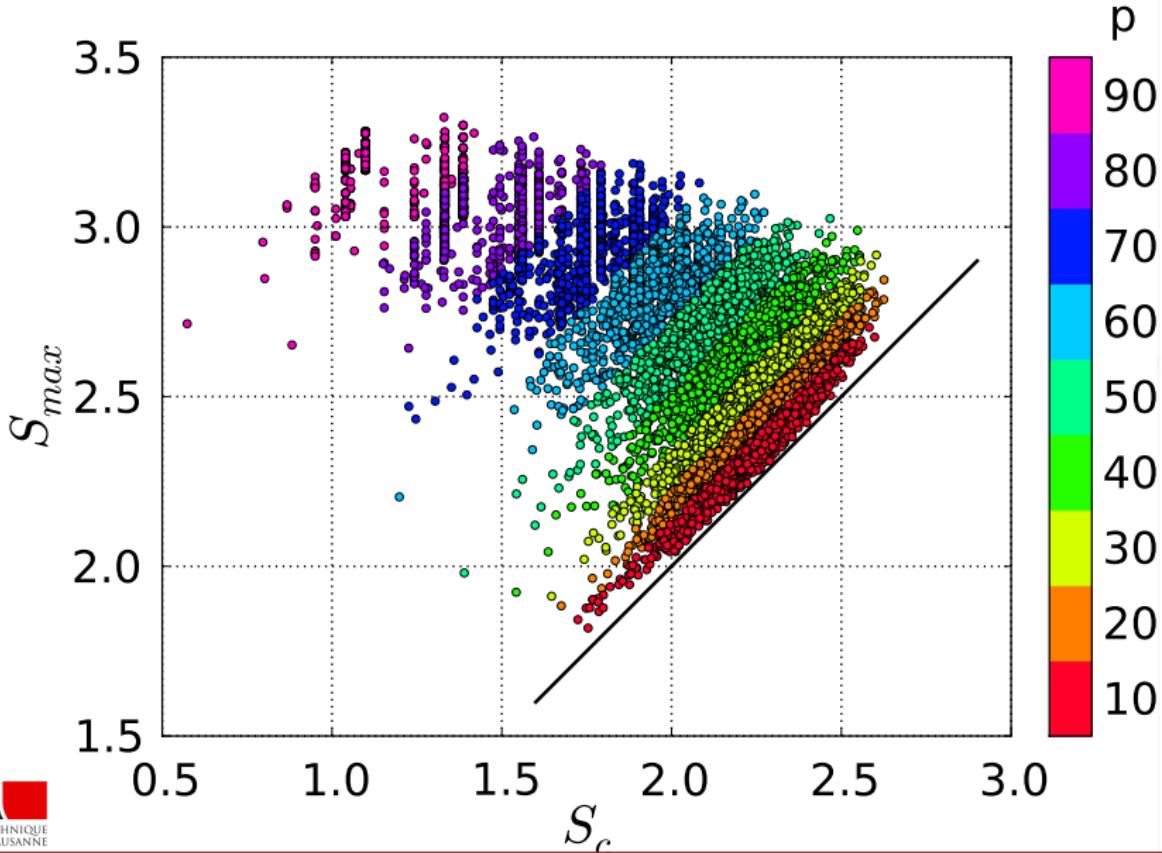
$$\sum_n p_n = 1$$

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

$$\sqrt{\sum_n p_n (\ln n - \langle \ln n \rangle)^2} = \sigma_{\ln n}$$

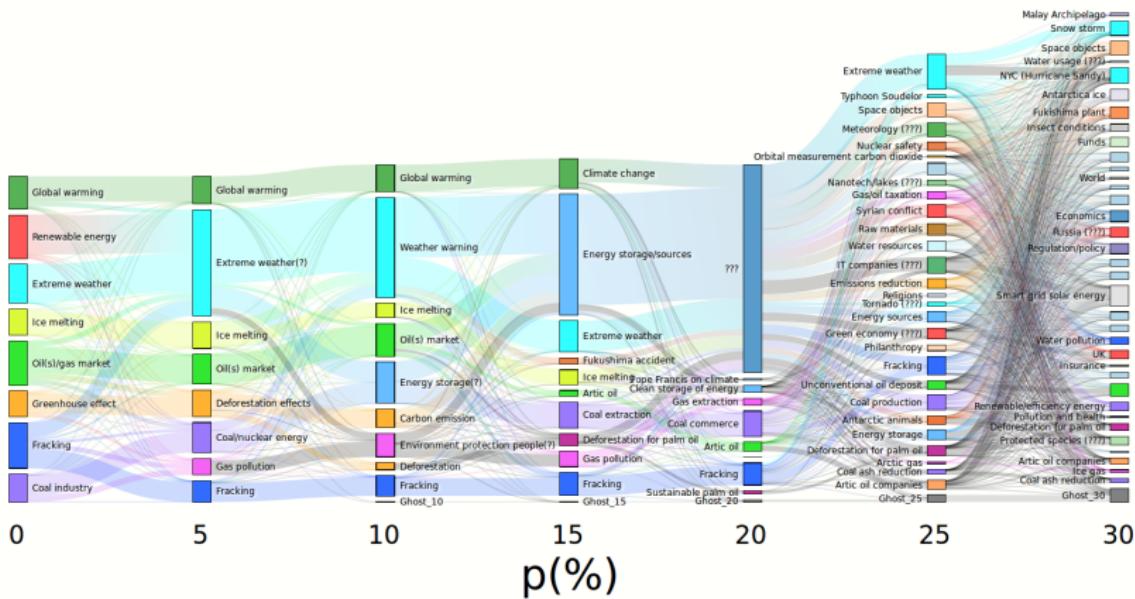


# Climate dataset



# Climate dataset

Data: Climate webdocs  $\max\min_{thr} = 0.005$   $w_{min} = 0.01$



# Optimal filtering

