

# Will This Paper Increase Your $h$ -index? Scientific Impact Prediction

Yuxiao Dong, Reid A. Johnson, Nitesh V. Chawla



Interdisciplinary Center for Network Science and Applications



# Scientific Impact



*Integral to the success of scientific research is the publication and dissemination of impactful work and findings.*



# Scientific Impact

*“An emerging area of interest in research  
on the ‘**science of science**’  
is the **prediction** of future **impact**.”*

*How?* *What?*

*J. A. Evans.  
Science 342, 2013*

- D. E. Acuna, S. Allesina, K. P. Kording. Future Impact: Predicting Scientific Success. Nature 489, 2012
- D. Wang, C. Song, A.-L. Barabasi. Quantifying long-term scientific impact. Science 342, 2013.
- B. Uzzi, S. Mukherjee, M. Stringre, B. Jones. Atypical Combinations and Scientific Impact. Science 342, 2013.
- H.-W. Shen and A.-L. Barabási. **Collective credit allocation in science**. PNAS 111, 2014.

# Academic Data



➤ A real-world academic dataset from



- 1,712,433 authors
- 2,092,356 papers
- 4,258,615 collaborations
- 8,024,869 citations
- <http://arnetminer.org/AMinerNetwork>

Arnetminer

MINING DEEP KNOWLEDGE FROM SCIENTIFIC NETWORKS

Whatever comes to your mind

Hot Topics

- Data Mining
- Machine Learning
- Social Network
- Deep Learning

Statistics

Researchers:	39,124,022
Publications:	79,029,784
Conferences/Journals:	330,236
Citations:	133,196,029
Knowledge Concepts:	7,854,301

# Scientific Impact: #citations

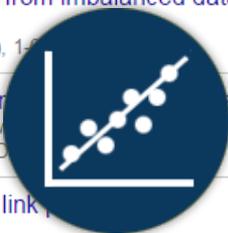
The number of citations of each publication

Title 1-20	Cited by	Year
<a href="#">SMOTE: synthetic minority over-sampling technique</a> NV Chawla, KW Bowyer, LO Hall, WP Kegelmeyer Journal of Artificial Intelligence Research (JAIR) 16, 321-357	2471 *	2002
<a href="#">Editorial: special issue on learning from imbalanced data sets</a> NV Chawla, N Japkowicz, A Kotcz ACM Sigkdd Explorations Newsletter 6 (1), 1-6	882	2004
<a href="#">SMOTEBoost: Improving prediction of the minority class in boosting</a> NV Chawla, A Lazarevic, LO Hall, KW Bowyer Knowledge Discovery in Databases: PKDD 2003, 107-119	450	2003
<a href="#">New perspectives and methods in link prediction</a> RN Lichtenwalter, JT Lussier, NV Chawla Proceedings of the 16th ACM SIGKDD international conference on Knowledge ...	229	2010
<a href="#">SVMs modeling for highly imbalanced classification</a> Y Tang, YQ Zhang, NV Chawla, S Krasser Systems, Man, and Cybernetics, Part B: Cybernetics, IEEE Transactions on 39 ...	218	2009

## #citations prediction

## Predicting the number of citations of publications

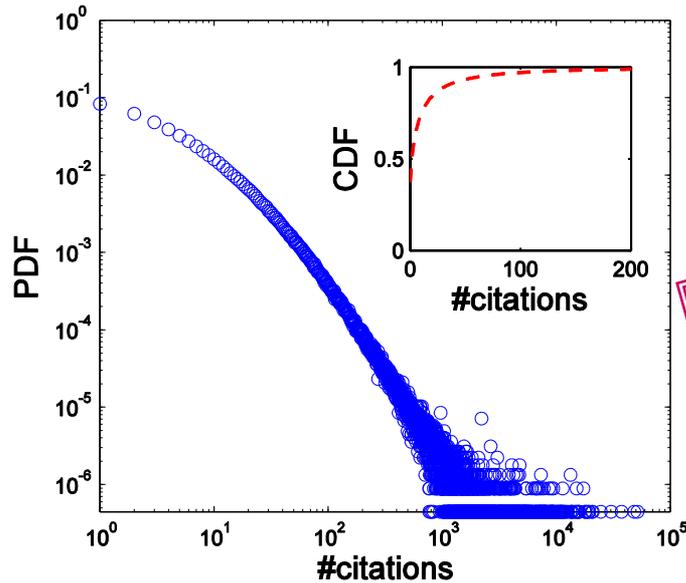
Title 1-20	Cited by	Year
<a href="#">SMOTE: synthetic minority over-sampling technique</a> NV Chawla, KW Bowyer, LO Hall, WP Kegelmeyer Journal of Artificial Intelligence Research (JAIR) 16, 321-357	2471 *	2002
<a href="#">Editorial: special issue on learning from imbalanced data sets</a> NV Chawla, N Japkowicz, A Kotcz ACM Sigkdd Explorations Newsletter 6 (1), 1-1	882	2004
<a href="#">SMOTEBoost: Improving prediction performance in boosting</a> NV Chawla, A Lazarevic, LO Hall, KW Bowyer Knowledge Discovery in Databases: PKDD	450	2003
<a href="#">New perspectives and methods in link prediction</a> RN Lichtenwalter, JT Lussier, NV Chawla Proceedings of the 16th ACM SIGKDD international conference on Knowledge ...	229	2010
<a href="#">SVMs modeling for highly imbalanced classification</a> Y Tang, YQ Zhang, NV Chawla, S Krasser Systems, Man, and Cybernetics, Part B: Cybernetics, IEEE Transactions on 39 ...	218	2009





# #citations prediction

- publications with few citations are extremely common
- publications with many citations are relatively rare



**Challenge**

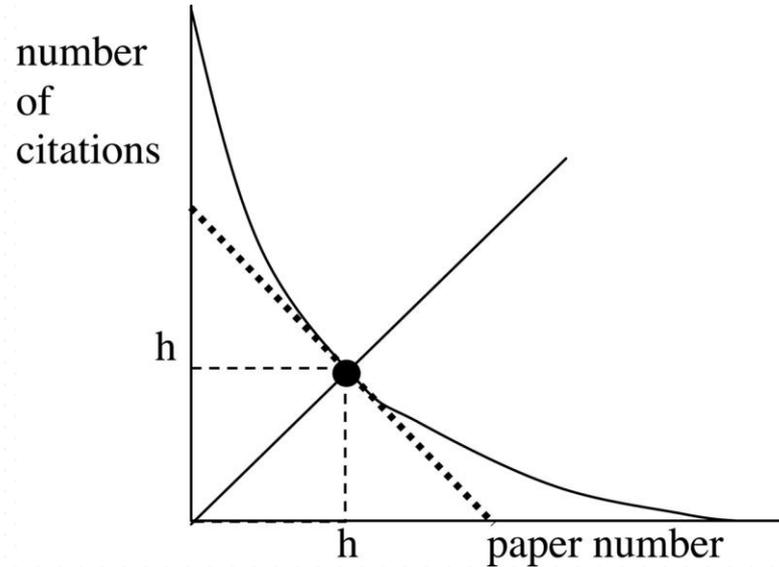


6.91% (155k out of 2 million) of the papers obtain more than 50 citations from 1950 to 2012.



# Scientific Impact: $h$ -index

## $h$ -index



# Scientific Impact: $h$ -index

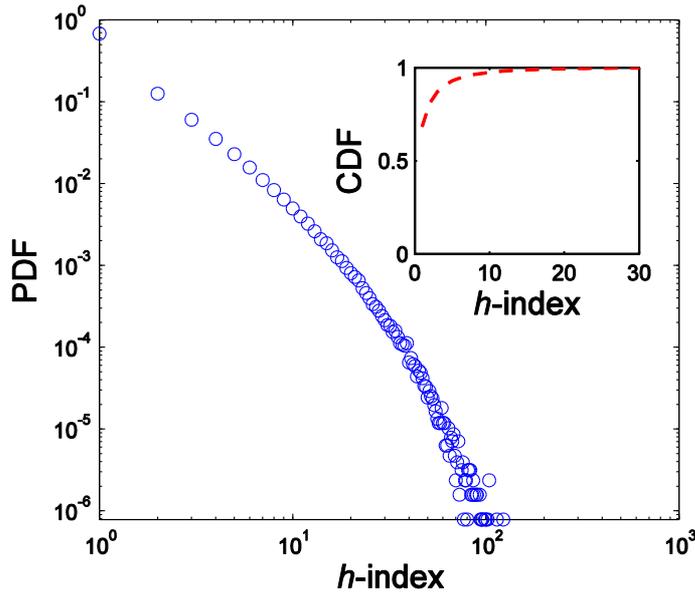
## The $h$ -index of each author

Experts	H-index	Rank
 <b>Thomas Huang</b> H-index: <b>128</b> , #Papers: <b>766</b> , #Citations: <b>65956</b> Professor, Beckman Institute at the University of Illinois	128	1
 <b>Anil K. Jain</b> H-index: <b>128</b> , #Papers: <b>440</b> , #Citations: <b>90064</b> Distinguished Professor, Michigan State University	128	2
 <b>Philip S. Yu</b> H-index: <b>127</b> , #Papers: <b>812</b> , #Citations: <b>68713</b> Professor and Wexler Chair in Information Technology, Department of Computer Science, University of Illinois Chicago	127	3
 <b>Jiawei Han</b> H-index: <b>116</b> , #Papers: <b>652</b> , #Citations: <b>90056</b> Professor, Department of Computer Science, University of Illinois at Urbana-Champaign	116	4
 <b>H. Garcia</b> H-index: <b>115</b> , #Papers: <b>429</b> , #Citations: <b>55531</b> Professor, Departments of Computer Science and Electrical Engineering, Stanford University	115	5



# $h$ -index prediction

Predicting the  $h$ -index of each author?



**Challenge**

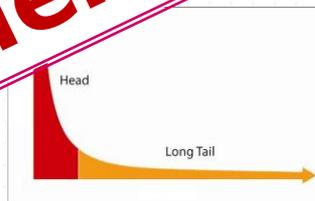


0.0125% (159 out of 1.7 million) of the researchers have an  $h$ -index over 60

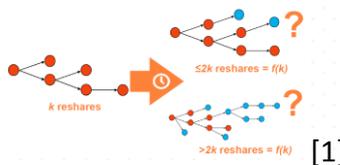
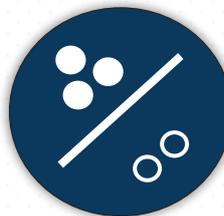


# Scientific Impact

- Predicting the #citations of each paper
- Predicting the *h*-index of each author



- Predicting whether a cascade will double in size<sup>[1]</sup>

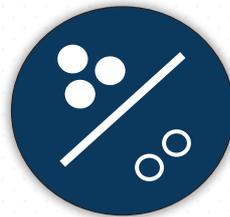


# Scientific Impact Prediction Problem



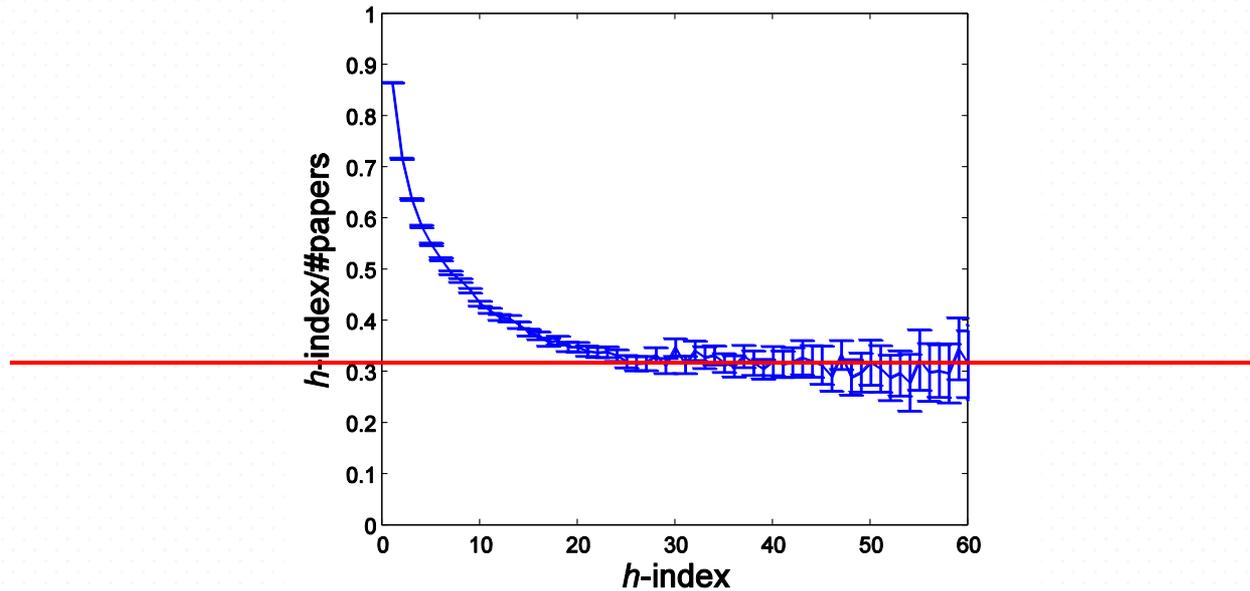
Given one paper and its author information,  
will it increase its **primary author's**  $h$ -index  
within a given time-frame  $\Delta t$ ?

the author of the given paper with the highest  $h$ -index.



# Scientific Impact: $h$ -index

## $h$ -index vs. $h$ -index/#papers



The ratio between one's  $h$ -index ( $\geq 20$ ) and her/his number of papers stabilizes at **0.3**.

# Scientific Impact Prediction Problem



## Can Cascades be Predicted?

Justin Cheng  
Stanford University  
jcccf@cs.stanford.edu

Lada A. Adamic  
Facebook  
ladamic@fb.com

P. Alex Dow  
Facebook  
adow@fb.com

primary author\*  
*h-index: 81*

Jon Kleinberg  
Cornell University  
kleinber@cs.cornell.edu

Jure Leskovec  
Stanford University  
jure@cs.stanford.edu

Given this paper at  $t=2014$  and its primary author, the task is to predict whether it will get at least 81 citations within  $\Delta t=5$  years.



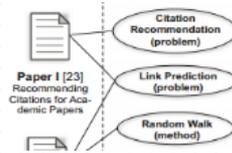
# Factors driving scientific impact



Author



Content



Reference



Temporal



Paper



Collaboration social network



Venue



# Factors --- author



**Author**

7 factors

## Can Cascades be Predicted?

first author	Justin Cheng Stanford University icccf@cs.stanford.edu	Lada A. Adamic Facebook ladamic@fb.com	P. Alex Dow Facebook adow@fb.com
primary author <i>h-index: 81</i>	Jon Kleinberg Cornell University kleinber@cs.cornell.edu	Jure Leskovec Stanford University jure@cs.stanford.edu	

all authors  
average author



# Factors --- content



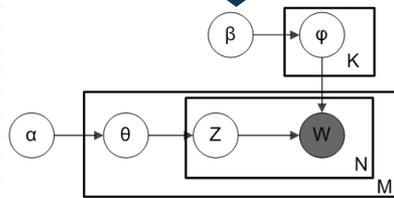
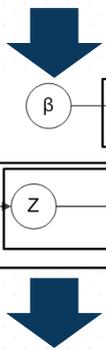
**Content**  
7 factors

## Will This Paper Increase Your h-index? Scientific Impact Prediction

Yuxiao Dong Dept. of Computer Science and Engineering, and CeNSA University of Notre Dame Notre Dame, IN 46556 ydong1@nd.edu  
 Reid A. Johnson Dept. of Computer Science and Engineering, and CeNSA University of Notre Dame Notre Dame, IN 46556 rjohns15@nd.edu  
 Nilesh V. Chawla Dept. of Computer Science and Engineering, and CeNSA University of Notre Dame Notre Dame, IN 46556 nchawla@nd.edu

**ABSTRACT**  
 Scientific impact plays a central role in the evaluation of the merit of scholars, departments, and institutions. A widely used measure of scientific impact is citations, with a growing body of literature focused on predicting the number of citations obtained by any given publication. The effectiveness of such predictions, however, is fundamentally limited by the power-law distribution of citations, whereby publications with few citations are extremely common and publications with many citations are relatively rare. Given this limitation, in this work we instead address a related question asked by many academic researchers in the course of writing a paper, namely, "Will this paper increase my h-index?" Using a real academic dataset with over 1.7 million authors, 2 million papers, and 8 million citation relationships from the premier online academic service ArXiv.org, we formulate a novel scientific impact prediction problem to examine several factors that can drive a paper to increase the primary author's h-index. We find that the researcher's authority on the publication topic and the venue in which the paper is published are crucial factors to the increase of the primary author's h-index, while the topic popularity and the co-authors' h-indices are of surprisingly little relevance. By leveraging relevant

**1. INTRODUCTION**  
 Integral to the success of scientific research is the publication and dissemination of impactful work and findings. Every scientific researcher strives for a high and desirable rank on an ever-expanding body of literature through a personal track-record of academic publications. The impact of such of these publications—both as a field of research and, by extension, the reputation of the author—can be influenced by a variety of factors. For example, significant research work may stem from a small number of exploratory papers that build up to pioneering work within a field, resulting in a series of less impactful papers that serve as stepping-stones to those of greater impact. Or a researcher may publish in different fields each with differing audience and levels of popularity, resulting in publications on some topics receiving more attention than those on others. Or, along with many publications that incrementally advance a field, a researcher may produce a groundbreaking work that transforms the field or even stimulates a new research area. As a result of such factors, a researcher's body of work is likely to be comprised of publications of varying impact. Accordingly, the impact of any particular publication can be difficult to predict.



scientific impact: 0.5  
 science of science: 0.4  
 social network: 0.1

topic popularity  
 deep learning is hot!

topic novelty

divergence of topics between this paper and its reference

topic diversity

divergence of topics of this paper

topic authority

authors' authority on the topics of this paper





# Factors --- venue



**Venue**  
2 factors

Top publications - Data Mining & Analysis

[Learn more](#)

Publication	h5-index
1. ACM SIGKDD International Conference on Knowledge discovery and data mining	69
2. IEEE Transactions on Knowledge and Data Engineering	57
3. ACM International Conference on Web Search and Data Mining	54
4. ACM Conference on Recommender Systems	36
5. IEEE International Conference on Data Mining (ICDM)	36
6. SIAM International Conference on Data Mining	35

average citations of papers in this venue

*h*-index contribution ratio of papers in this venue

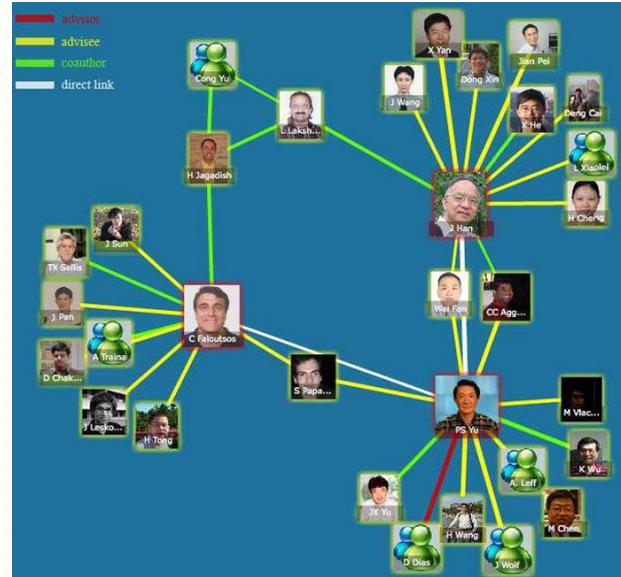




# Factors --- social

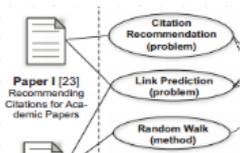


Collaboration social network  
4 factors



degree  
Pagerank  
coauthors' h-indices

# Factors --- reference



## Reference

2 factors

### 7. REFERENCES

- [1] M. Ahmed, S. Spagna, F. Huici, and S. Niccolini. A peek into the future: Predicting the evolution of popularity in user generated content. In *WSDM '13*, pages 607–616. ACM, 2013.
- [2] S. Bethard and D. Jurafsky. Who should I cite: Learning literature search models from citation behavior. In *CIKM '10*, pages 609–618. ACM, 2010.
- [3] D. M. Blei, A. Y. Ng, and M. I. Jordan. Latent Dirichlet allocation. *JMLR*, 3:993–1022, 2003.
- [4] C. Castillo, D. Donato, and A. Gionis. Estimating the number of citations using author reputation. In *SPIRE '07*, pages 107–117. Springer, 2007.
- [5] J. Cheng, L. Adamic, P. A. Dow, J. M. Kleinberg, and J. Leskovec. Can cascades be predicted? In *WWW '14*, pages 925–936, 2014.
- [6] E. Garfield. Citation indexes for science: A new dimension in documentation through association of ideas. *Science*, 122(3159):108–111, 1955.

citations of references

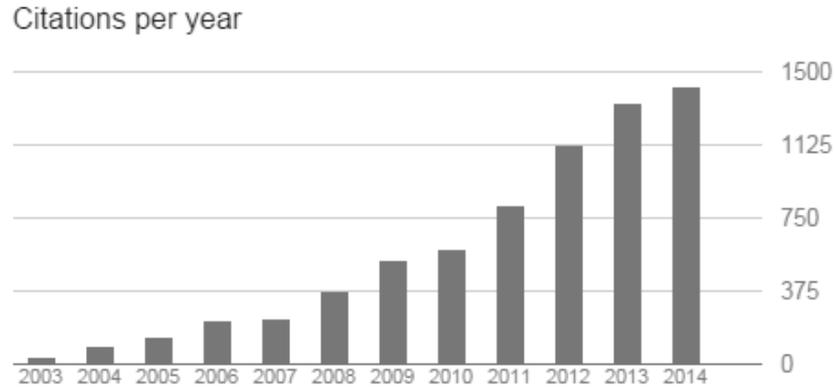
*h*-index of references



# Factors --- temporal



**Temporal**  
4 factors



authors' *h*-index increasing rate

# Factor Definition

**Table 1: Factor Definition.** We employ six categories of factors, comprised of author, topic, reference, social, venue, and temporal attributes. *max-h-index* denotes the *h-index* of the primary author (i.e., the author with the maximum *h-index*) of a given paper.

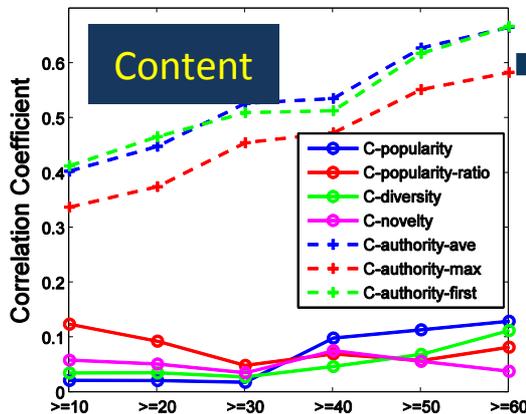
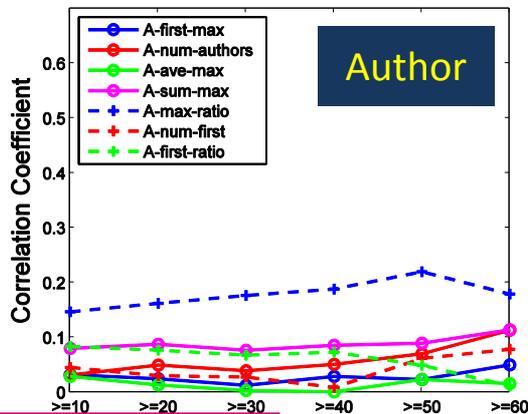
	Factor	Description
Author	<i>A-first-max</i>	The first author's <i>h-index</i> divided by the <i>max-h-index</i> .
	<i>A-ave-max</i>	The average <i>h-index</i> of all authors divided by the <i>max-h-index</i> .
	<i>A-sum-max</i>	The sum of <i>h-indices</i> divided by the <i>max-h-index</i> .
	<i>A-first-ratio</i>	The ratio between <i>max-h-index</i> and the number of papers attributed to the first author.
	<i>A-max-ratio</i>	The ratio between <i>max-h-index</i> and the number of papers attributed to the primary author.
	<i>A-num-authors</i>	The number of authors of the given paper.
	<i>A-num-first</i>	The number of papers by the first author.
Content	<i>C-popularity</i>	The average number of citations over different topics (see Eq. 1).
	<i>C-popularity-ratio</i>	The average number of citations over different topics divided by the <i>max-h-index</i> .
	<i>C-novelty</i>	The topic novelty of this paper (see Eq. 2).
	<i>C-diversity</i>	The topic diversity of this paper (see Eq. 3).
	<i>C-authority-first</i>	The consistence between the first author's authority and this paper (see Eq. 4).
	<i>C-authority-max</i>	The consistence between the primary author's authority and this paper.
	<i>C-authority-ave</i>	The average consistence between each author's authority and this paper.
Venue	<i>V-ratio-max</i>	The ratio between the number of papers $\geq$ <i>max-h-index</i> citations divided by the <i>max-h-index</i> .
	<i>V-citation</i>	The average number of citations of all references divided by the <i>max-h-index</i> .
Social	<i>S-degree</i>	The number of co-authors of the paper's authors.
	<i>S-pagerank</i>	The PageRank values of the paper's authors in the weighted collaboration network.
	<i>S-h-co-author</i>	The average <i>h-index</i> of co-authors of the paper's authors divided by the <i>max-h-index</i> .
	<i>S-h-weight</i>	The weighted average <i>h-index</i> of co-authors of the paper's authors divided by the <i>max-h-index</i> .
Reference	<i>R-ratio-max</i>	The ratio between the number of references $\geq$ <i>max-h-index</i> and the total number of references.
	<i>R-citation</i>	The average number of citations divided by the maximum <i>h-index</i> .
Temporal	<i>T-ave-h</i>	The average $\Delta h$ -indices of the authors between now and three years ago.
	<i>T-max-h</i>	The maximum $\Delta h$ -index between now and three years ago.
	<i>T-h-first</i>	The $\Delta h$ -index of the first author between now and three years ago.
	<i>T-h-max</i>	The $\Delta h$ -index of the <i>max-h-index</i> author between now and three years ago.

6 groups

26 factors

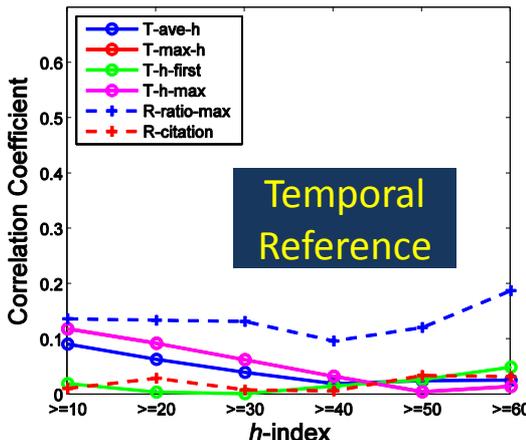
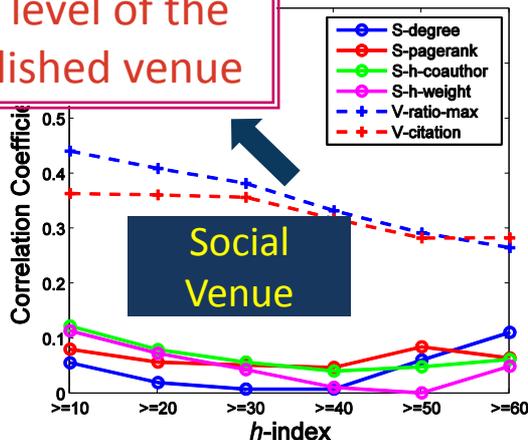


# Factors Correlation



authors' authority on the topics of this paper

$t = 2007$   
 $\Delta t = 5$



X-axis:  
primary author's  $h$ -index

Y-axis:  
correlation coefficient

the level of the published venue





# Factors Correlation

A scientific **researcher's authority** on a topic is the most decisive factor in facilitating an increase in his or her *h*-index.

# Factors Correlation



The level of **the venue** in which a given paper is published is another crucial factor in determining the probability that it will contribute to its authors'  $h$ -indices.



# Factors Correlation

Publishing on an academically “hot” but unfamiliar topic is difficult to further one's scientific impact, at least as measured by an increase in one's  $h$ -index.

# Prediction: predictability



Is Scientific Impact Predictable?

# Prediction: predictability

$t = 2007$      $\Delta t = 5$   
21,519 papers

On average, 30.5% of papers successfully contributed to their primary author's  $h$ -indices in 2012.

Task: predict whether the number of citations for each paper published in 2007 is larger than or equal to the primary author's  $h$ -index in 2012

R: Random guess

Features: 26 factors

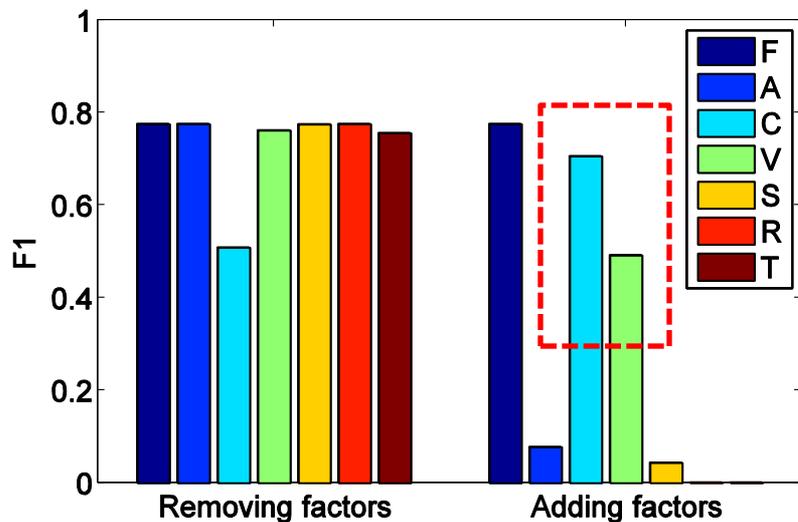
LRC: Logistic regression

Half training, half test

Method	Pre.	Rec.	$F_1$	AUC	Accu.	Pre@3	MAP
R	0.305	0.500	0.375	0.500	0.500	0.674	0.522
LRC	0.854	0.711	0.776	0.938	0.875	0.925	0.965



# Prediction: factor contribution



- F: Full factors
- A: Author
- C: Content
- V: Venue
- S: Social
- R: Reference
- T: Temporal

$t = 2007$

$\Delta t = 5$

Logistic regression



# Prediction: predictability

## Can Cascades be Predicted?

Justin Cheng  
Stanford University  
jcccf@cs.stanford.edu

Lada A. Adamic  
Facebook  
ladamic@fb.com

P. Alex Dow  
Facebook  
adow@fb.com

Jon Kleinberg  
Cornell University  
kleinber@cs.cornell.edu

Jure Leskovec  
Stanford University  
jure@cs.stanford.edu

Published at 2014

$\Delta t = 5$  years

$\Delta t = 10$  years

Is a paper more predictable  
given a long or short timeframe  $\Delta t$ ?



# Prediction: predictability

Published at 2014

## Inferring User Demographics and Social Strategies in Mobile Social Networks

Yuxiao Dong<sup>1</sup>, Yang Yang<sup>2</sup>, Jie Tang<sup>3</sup>, Yang Yang<sup>1</sup>, Nitesh V. Chawla<sup>1</sup>

<sup>1</sup>Department of Computer Science and Engineering, University of Notre Dame

<sup>2</sup>Interdisciplinary Center for Network Science and Applications, University of Notre Dame

<sup>3</sup>Department of Computer Science and Technology, Tsinghua University

ydong1@nd.edu, yyang.thu@gmail.com, jetang@tsinghua.edu.cn, yyang1@nd.edu, nchawla@nd.edu

## Can Cascades be Predicted?

Justin Cheng  
Stanford University  
jcccf@cs.stanford.edu

Lada A. Adamic  
Facebook  
ladamic@fb.com

P. Alex Dow  
Facebook  
adow@fb.com

Jon Kleinberg  
Cornell University  
kleinber@cs.cornell.edu

Jure Leskovec  
Stanford University  
jure@cs.stanford.edu

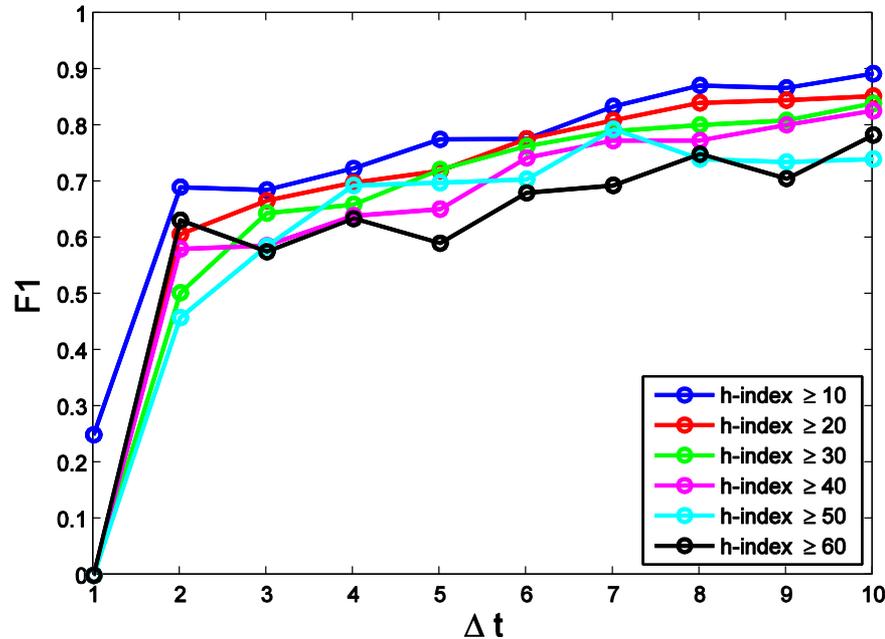
Primary author's  $h$ -index: 33

Primary author's  $h$ -index: 81

Is a primary author with a high or  
a low  $h$ -index more predictable?



# Prediction: predictability



$t + \Delta t = 2012$   
Logistic regression

1. more difficult for papers with a high  $h$ -index primary author
2. more difficult when given a shorter timeframe  $\Delta t$ .



# Future work

1. *Only work on computer science domain*

*TODO: physics, mathematics, biology ...*

2. *Authors' h-indices evolve within  $\Delta t$*

*TODO: co-evolution of authors' h-indices and #citations*



*When a measure becomes a target, it ceases to be a good measure*

*---Charles Goodhart*

# Acknowledgements

*Army Research Laboratory  
(ARL)*

*U.S. Air Force Office of Scientific Research  
(AFOSR)*

*Defense Advanced Research Projects Agency  
(DARPA)*

*National Science Foundation  
(NSF)*



# Thanks

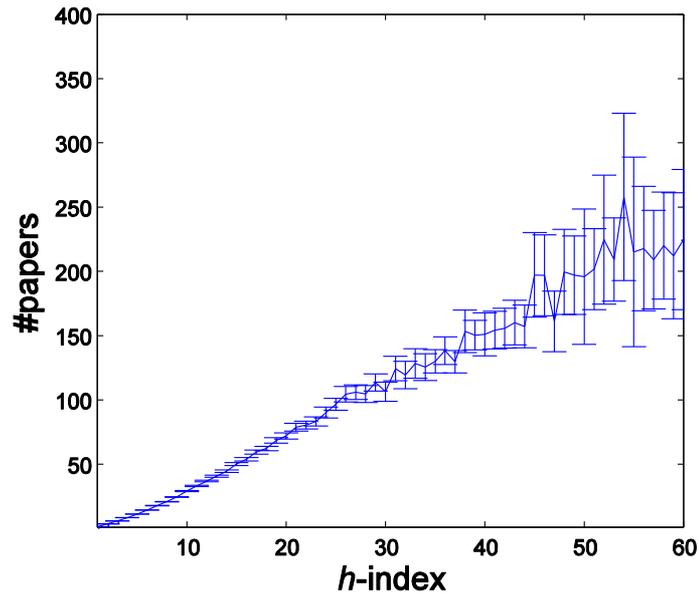
*Standing on the shoulders of giants*  
--- Isaac Newton

Q & A



# Scientific Impact: $h$ -index

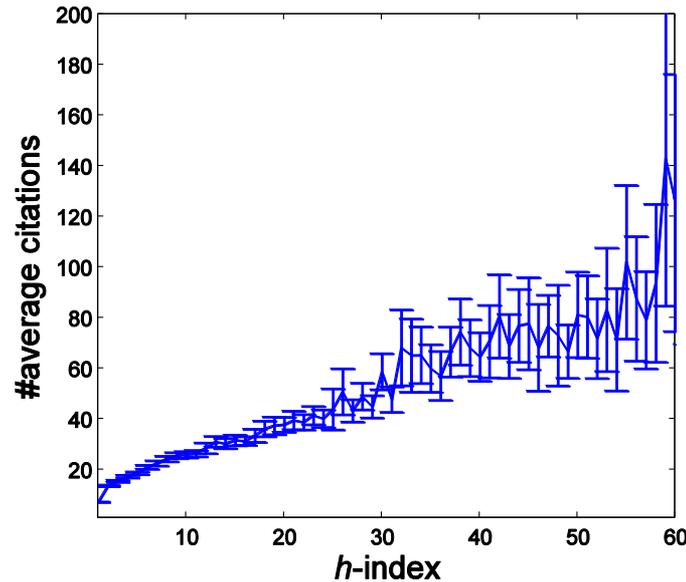
## $h$ -index vs. #papers





# Scientific Impact: $h$ -index

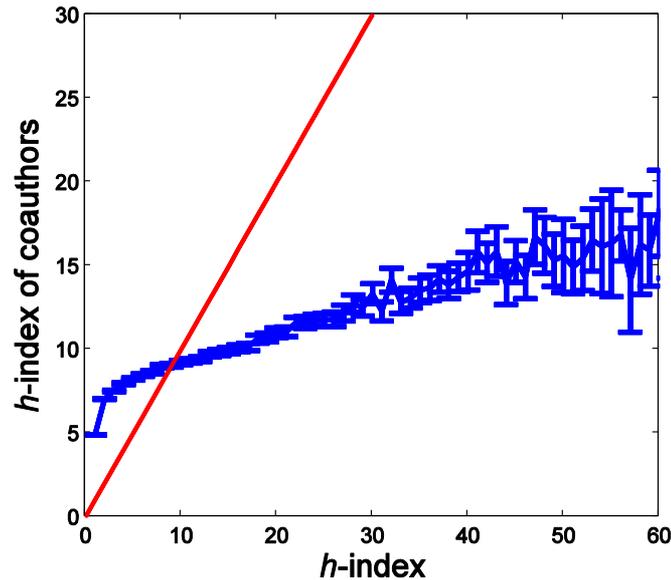
## $h$ -index vs. #average-citations



The average number of citations for each author is larger than her/his  $h$ -index.

# Scientific Impact: $h$ -index

## $h$ -index vs. average $h$ -index of coauthors

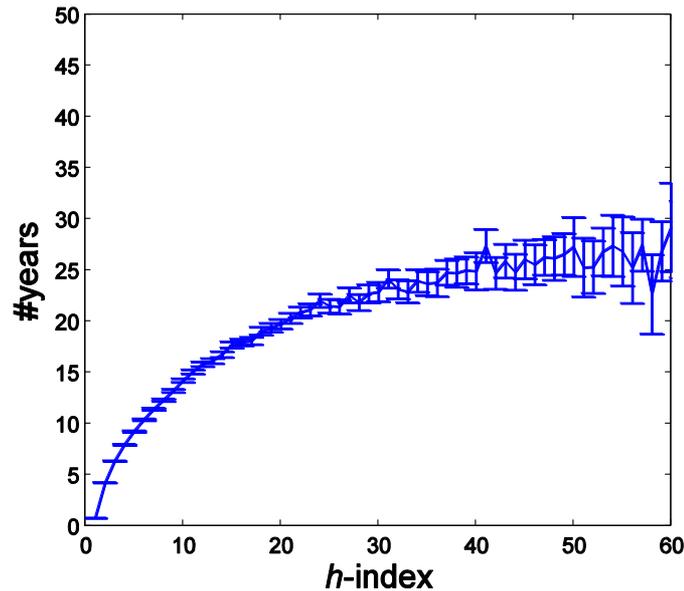


Typically, the author's  $h$ -index becomes larger than the co-authors'  $h$ -indices at the expected point of the author's Ph.D. graduation.



# Scientific Impact: $h$ -index

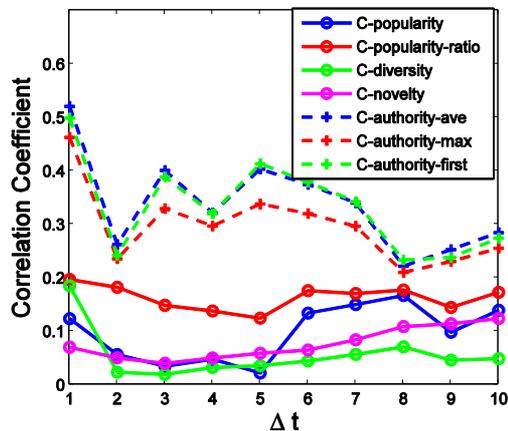
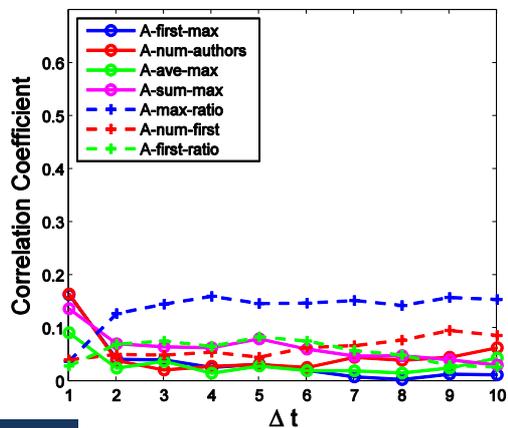
## $h$ -index vs. #career years



The rate at which the  $h$ -index increases itself increases as the length of time spent in academia becomes longer (i.e., the rich get richer).



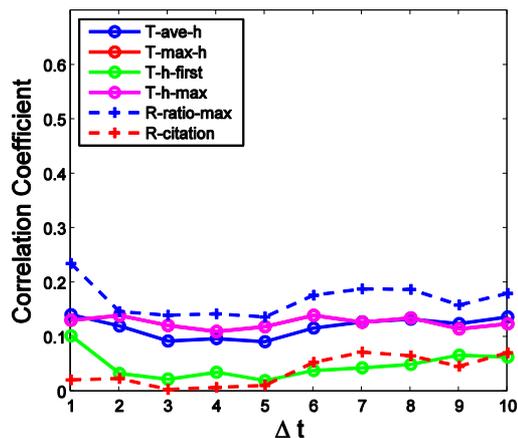
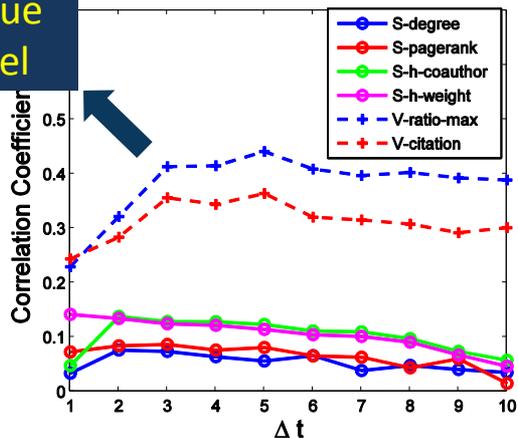
# Factors Correlation 2



authors' authority on the topics of this paper

t = 2002

venue level



X-axis:

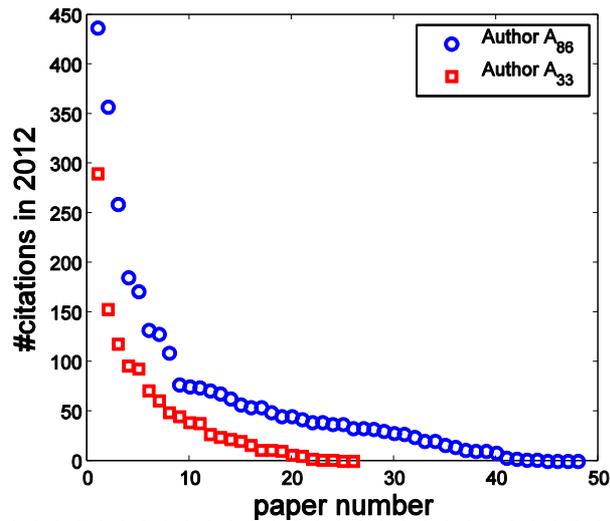
$\Delta t$

Y-axis:

correlation coefficient



# Prediction: case study 1



Two anonymous authors  
 $A_{86}$  and  $A_{33}$

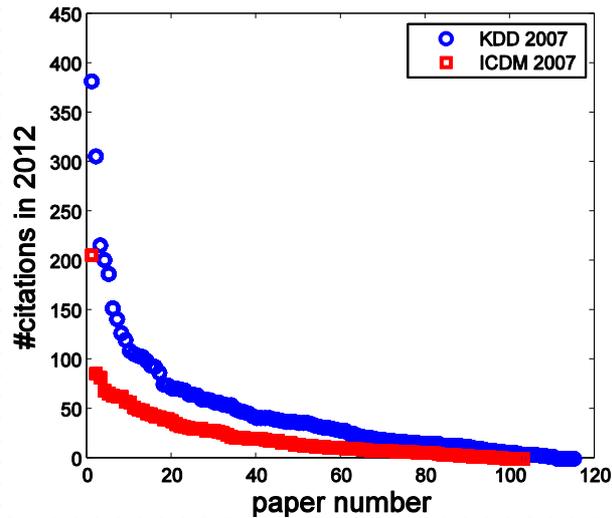
Authors	Pre.	Rec.	F <sub>1</sub>	AUC	Accu.	Pre@k	MAP
$A_{86}$	0.500	0.375	0.429	0.584	0.833	0.375	0.346
$A_{33}$	1.000	0.667	0.800	0.856	0.885	0.667	0.849

$t = 2007$

$\Delta t = 5$

Logistic regression

# Prediction: case study 2



*Two venues  
KDD and ICDM*

Venues	Pre.	Rec.	$F_1$	AUC	Accu.
KDD'07	0.800	0.889	0.842	0.884	0.818
ICDM'07	0.842	0.593	0.696	0.886	0.825

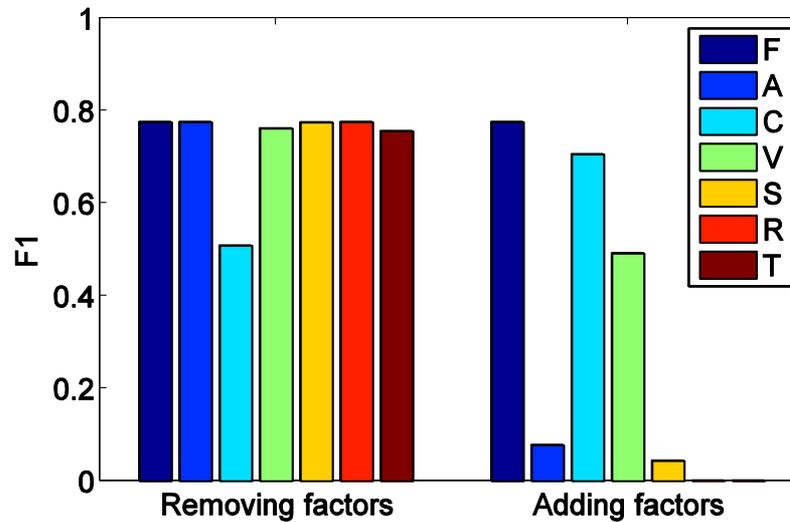
$t = 2007$

$\Delta t = 5$

Logistic regression



# Prediction: factor contribution



F: Full  
A: Author  
C: Content  
V: Venue  
S: Social  
R: Reference  
T: Temporal

$t = 2007$

$\Delta t = 5$

Logistic regression