

# Impact of lexical and sentiment factors on the popularity of scientific papers

**JULIAN SIENKIEWICZ\*** & **EDUARDO G. ALTMANN†**

**Max Planck Institute for the Physics of Complex Systems, Dresden**

\* current address: Faculty of Physics, Warsaw University of Technology

† current address: School of Mathematics and Statistics, University of Sydney

9 November 2016

## MOTIVATION

- 1 Citations for an article can be regarded as a **proxy for the attention** (popularity) the paper achieved in the scientific community.
- 2 Investigation how **textual properties** of scientific papers affect the adoption and spread of scientific results as quantified by the number of citations.
- 3 Last but not least: dealing with recent results (e.g., Letchford et al. R Soc Open Sci 2, 150266) showing that is the **negative correlation** between title length and citations (i.e., shorter the titles more citations).



## GOALS

### KEY FACTORS

Systematic investigation of how different **textual properties** of scientific papers such as

- text length,
- text complexity,
- sentiment

affect the number of **citations** they acquire.

In this way we want identify **key factors** that influence scientific popularity.

## GOALS

### KEY FACTORS

Systematic investigation of how different **textual properties** of scientific papers such as

- text length,
- text complexity,
- sentiment

affect the number of **citations** they acquire.

In this way we want identify **key factors** that influence scientific popularity.

### CITATION PATTERN DIFFERENCES

We want to show **differences** between **most cited** (top) and **typical** papers using **quantile regression** approach.

## DATA

### Web of Science service

#### FILTERING

- papers marked as articles published in the period of **1995—2004**
- papers needed to fulfil two conditions:
  - 1 journal active in **all** mentioned years (e.g., PLOS journals absent)
  - 2 there had to be at least **1.000 articles** in the given period (e.g., Rev Mod Phys absent)

## DATA

### Web of Science service

#### FILTERING

- papers marked as articles published in the period of **1995—2004**
- papers needed to fulfil two conditions:
  - 1 journal active in **all** mentioned years (e.g., PLOS journals absent)
  - 2 there had to be at least **1.000 articles** in the given period (e.g., Rev Mod Phys absent)

#### OUTCOME

- over **4.300.000** articles from over **1.500** different journals,
- information about the **title of the paper**, the number of its **authors**, full **abstract contents** and OECD category it had been classified to,
- the **number of citations** it acquired between being published and 31<sup>st</sup> December 2014

## MEASURES

property	title	abstract
length	number of characters	number of words
complexity	— z-index Herdan's $C$	Gunning fog index $F$ z-index Herdan's $C$
sentiment	valence arousal	valence arousal
number of authors		

① Fog index:  $F = \left( \frac{\#words}{\#sentences} + 100 \frac{\#complexwords}{\#words} \right)$

② Herdan's  $C$ :  $C = \frac{\log N}{\log M} \left[ \begin{array}{l} M - \text{text length} \\ N - \text{vocabulary size} \end{array} \right]$

③ z-index:  $z_{M,N} = \frac{N - \mu(M)}{\sigma(M)}$

④ Valence — emotional sign of the text (positive - 9, neutral - 5, negative - 1)

⑤ Arousal — level of emotional activation (low - 1, medium - 5, high - 9)



## QUANTILE REGRESSION (QR)

### IDEA

Find coefficients  $\alpha$  and  $\beta$  of the relation

$$Y = \alpha(\tau) + \beta(\tau)X$$

dividing the dataset so that  $\tau$  points lay below the line and  $(1 - \tau)$  are above it.

### ADVANTAGES

- we can examine different regimes (ranges) of  $Y$ ,
- the log of  $p$ -th quantile is equal to the  $p$ -th quantile log-transformed  $Y$

# QUANTILE REGRESSION (QR)

## IDEA

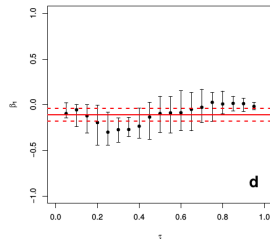
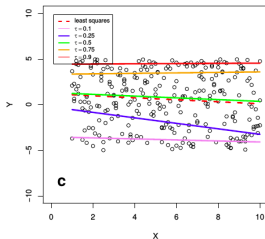
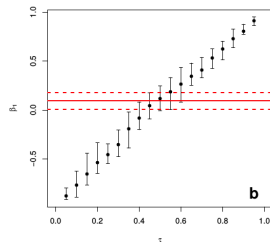
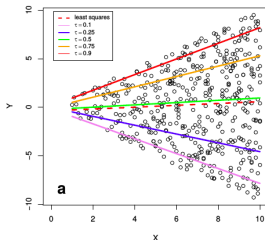
Find coefficients  $\alpha$  and  $\beta$  of the relation

$$Y = \alpha(\tau) + \beta(\tau)X$$

dividing the dataset so that  $\tau$  points lay below the line and  $(1 - \tau)$  are above it.

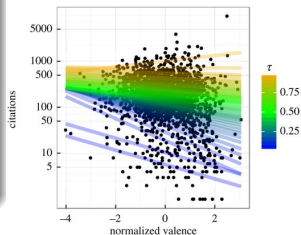
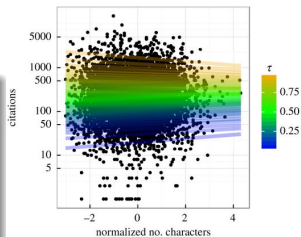
## ADVANTAGES

- we can examine different regimes (ranges) of  $Y$ ,
- the log of  $p$ -th quantile is equal to the  $p$ -th quantile log-transformed  $Y$

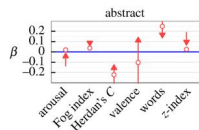
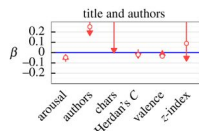
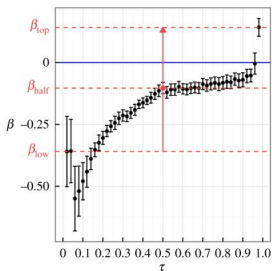
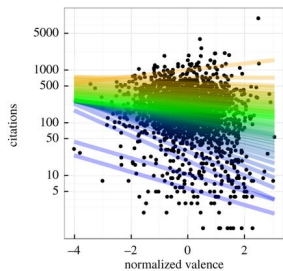
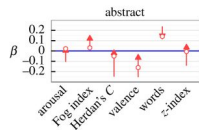
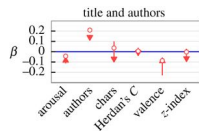
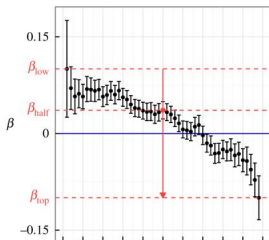
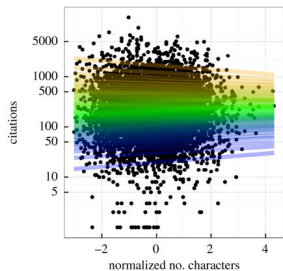


## RESULTS - QR - DISCUSSION

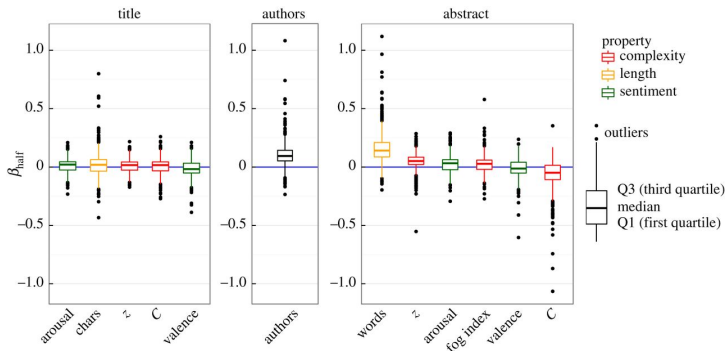
- broad scattering of the points — visual inspection fails even to detect whether the relation between  $X$  and  $Y$  is positive or negative,
- Pearson correlation coefficient  $r$  yields:  
 $r = 0.02 \pm 0.01$  for title length (Science) and  
 $r = -0.21 \pm 0.03$  for valence (Nature Genetics),
- such difference motivates us to go beyond linear correlations, which rely on the (homoscedasticity) assumption of uniform errors in the whole dataset.



## RESULTS - QR



## RESULTS - COMPARISON OF FACTORS



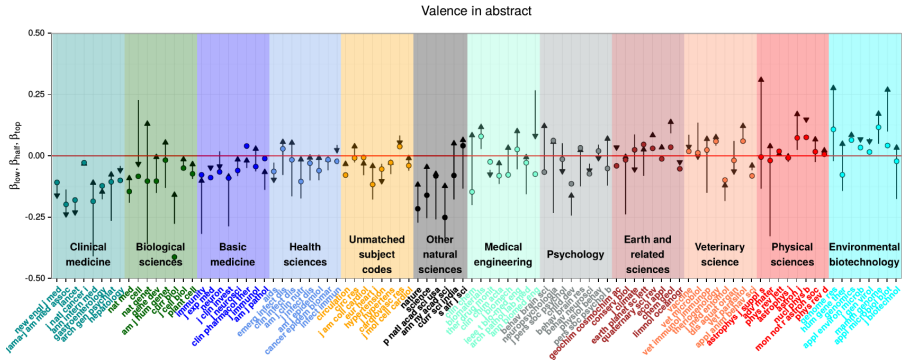
- The influence of factors is overall rather weak -  $|\beta| < 0.5$  ( $\beta = \ln 2$  means that the number of citations  $Y$  doubles by moving 1 standard deviation in  $X$ ).
- the strongest factors are (i) the number of words in the abstract, (ii) the number of authors, and (iii) z-index in the abstract (over 75% of journals — equivalently, the whole box, are placed above zero).
- factors in the abstract are more visible than in the title

## RESULTS - DIFFERENCE BETWEEN TYPICAL AND TOP PAPERS

property	factor	$\beta_{top} > \beta_{half}$	$\beta_{top} < \beta_{half}$	$\beta_{top} \neq \beta_{half}$
length	no. of characters (title)	2.6%	44.4%	47.0%
	no. of words (abstract)	8.3%	29.4%	36.7%
	mean			41.9%
complexity	Herdan's $C$ (title)	18.7%	8.5%	27.2%
	Herdan's $C$ (abstract)	34.9%	6.5%	41.4%
	z-index (title)	8.3%	16.7%	25.0%
	z-index (abstract)	24.6%	7.7%	32.3%
	fog index (abstract)	26.4%	8.0%	34.4%
	mean			32.0%
sentiment	arousal (title)	11.0%	13.5%	24.5%
	arousal (abstract)	15.7%	13.7%	29.4%
	valence (title)	16.1%	11.3%	27.4%
	valence (abstract)	29.2%	5.7%	34.9%
	mean			29.1%
	no. of authors	4.0%	39.6%	43.6%
	overall mean			33.7%



## RESULTS - COMPARISON ACROSS JOURNALS (VALENCE IN ABSTRACT)



Variation across journals is partially explained by disciplines, e.g. for *clinical medicine* all values of  $\beta$  in the case of valence in abstract are below zero, whereas for *physical sciences*, the majority is positive.



## SUMMARY

- 1 Investigation on how textual properties of scientific papers relate to the number of citations they receive,
- 2 Main finding: correlations are **non-linear** and affect differently **most-cited** and **typical** papers,
- 3 In most journals **short titles** correlate **positively** with citations only for the **most** cited papers, for typical papers the correlation is in most cases **negative**,
- 4 Statistically significant effect present for **most factors**, but it is typically **weak** ( $|\beta| < 0.5$ ),
- 5 large **variability** across journals

details & some data: **R Soc Open Sci 3, 160140 (2016)**