Sic Transit Gloria Manuscriptum: Two Views of the Aggregate Fate of Ancient Papers

Mayank Singh¹, Rajdeep Sarkar², Pawan Goyal¹, Animesh Mukherjee¹ and Soumen Chakrabarti³

^{1,3}Department of Computer Science and Engineering

²Department of Mathematics

^{1,2}Indian Institute of Technology, Kharagpur, India

³Indian Institute of Technology, Bombay, India

mayank.singh@cse.iitkgp.ernet.in, rajdeep.sarkar@iitkgp.ac.in

{pawang,animeshm}@cse.iitkgp.ernet.in, soumen.chakrabarti@gmail.com

July 23, 2018

Abstract

When PageRank began to be used for ranking in Web search, a concern soon arose that older pages have an inherent — and potentially unfair — advantage over emerging pages of high quality, because they have had more time to acquire hyperlink citations. Algorithms were then proposed to compensate for this effect. Curiously, in bibliometry, the opposite concern has often been raised: that a growing body of recent papers crowds out older papers, resulting in a collective amnesia in research communities, which potentially leads to reinventions, redundancies, and missed opportunities to connect ideas. A recent paper by Verstak et al. reported experiments on Google Scholar data, which seemed to refute the amnesia, or aging, hypothesis. They claimed that more recently written papers have a larger fraction of outbound citations targeting papers that are older by a fixed number of years, indicating that ancient papers are alive and well-loved and increasingly easily found, thanks in part to Google Scholar. In this paper we show that the full picture is considerably more nuanced. Specifically, the fate of a fixed sample of papers, as they age, is rather different from what Verstak et al.'s study suggests: there is clear and steady abandonment in favor of citations to newer papers. The two apparently contradictory views are reconciled by the realization that, as time passes, the number of papers older than a fixed number of years grows rapidly.

1 Introduction

The volume of scholarly publication per year has grown dramatically in the last two decades, along with an explosion in the number of papers that are available online, indexed, and searchable. The effect of such enhanced access on subsequent research is not entirely clear: there have been recent conflicting claims. Parolo et al. [3] present evidence that it is becoming "increasingly difficult for researchers to keep track of all the publications relevant to their work". Based on analysis of citation data, they propose a pattern of a paper's citation counts per year, which peaks within a few years and then the typical paper fades into obscurity. This work has seen considerable press following, with headlines ranging from the tongue-in-cheek "Study shows there are too many studies" to the more alarmist "Science is 'in decay' because there are too many studies". Chakraborty et al. [2] present a more nuanced analysis that clusters papers into the ephemeral and the enduring, giving some hope that not all creativity is lost in the sands of time. Meanwhile, Verstak et al. [5], from the Google Scholar team, claim that fear of evanescence is

misplaced, and that older papers account for an increasing fraction of citations as time passes. Superficially, these claims seem to be at odds with each other. But their simultaneous validity is easily explained by mild additional probing of citation statistics, which we present in this paper. Our main observations are as follows.

- We confirm that the fraction of citations made in a paper p_0 to other papers that are older than p_0 by a fixed time window (say, 10 to 15 years older) grows as we sample p_0 from more and more recent years.
- But this happens because the number of papers becoming older than a threshold is steeply increasing with time.
- If instead we fix a (sampled) set of papers and trace the rate of citations to them as time progresses, aging and passing into obscurity are strongly visible.
- These observations for citing behavior hold true across many disciplines (computer science, biomedical), sub-disciplines (algorithms, theory) as well as conferences.

2 Dataset

We begin our analysis using a dataset from the computer science domain, crawled from Microsoft Academic Search (MAS)¹. The dataset consists of bibliographic information of papers, the title of the paper, a unique index for the paper, its author(s), the affiliation of the author(s), the year of publication, the publication venue, references, citation contexts, the related field(s) of the paper, the abstract and the keywords of the papers [1]. In addition, we also use another dataset available from the biomedical domain². Some general statistics of both data sets are shown in Table 1.

Table 1: General statistics about the Computer Science (CS) and Biomedical (BM) dataset.

	CS	BM
Sub-fields	25	1
Publication count	1,359,338	801,252
Author count	138,923	1,985,890
Year range	1960-2010	1996-2014

3 Analysis of Citation Behavior

We analyze the data sets in various ways to investigate the claim by Verstak et al. [5]. First of all, we show that if we divide the cited papers into two time-zones – those having a time difference of $\leq t$ years with the citer paper and all others, we obtain similar results as described in [5]. However, if we fix the set of old papers, our results suggest that the number of citations that these papers get over the years are indeed affected by the aging behavior.

3.1 Fraction of citations to (all) 'old' papers

For the papers published between 1970-2010, Figure 1(a) shows fraction of out-citations given to all the papers older than t years for three different values of t (10, 15, 20). For a particular year y and time-window t, we compute fraction of older citations $R_{t,y}$ as

$$R_{t,y} = \frac{C_t}{C_y} \tag{1}$$

where C_y represents number of citations from papers published in year y and C_t represents number of citations to papers with publication time difference $\geq t$ from the year y. We observe that this is consistent with the claim by Verstak et al. [5], fraction of citations to the older papers is increasing over the years for all the values of t.

¹http://academic.research.microsoft.com

²http://www.ncbi.nlm.nih.gov/pmc/tools/ftp, downloaded in May, 2014

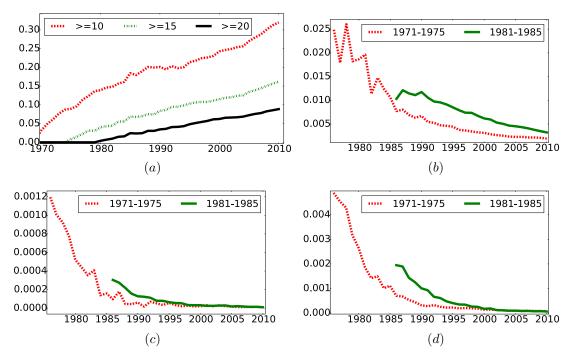


Figure 1: Fraction of citations between (a) 1970-2010 to all the older papers in three different time-windows (\geq 10 years, \geq 15 years and \geq 20 years, (b) 1975-2010 to a fixed set of 100 most cited papers from 1971-1975 and 1981-1985, (c) 1975-2010 to a fixed set of 100 random papers from 1971-1975 and 1981-1985 and (d) 1975-2010 to a fixed set of 500 random papers from 1971-1975 and 1981-1985.

Our central claim is that while the fraction of citations to all the older papers might be increasing over time, the set of older papers which are getting more citations is not the same. The above formulation fails to capture the notion that in each successive year, an increasing number of publication set is added to time-window $\geq t$. Thus, the increase in the fraction of citations might just be a byproduct of the increase in the number of publications. We suggest that if instead, the number of publications are fixed for various time-windows and we observe the fraction of citations going to these papers, the results give a very different picture.

As a first experiment, we took 100 most cited articles from the years 1971-1975 and 1981-1985 each. Figure 1(b) presents fraction of citations going to each of these sets over the years (i.e., what fraction of the out-going citations in that year went to these set of papers). We clearly see that this fraction is decreasing over the time. Further, we conducted experiments for random set of papers for each year intervals. Figure 1(c and d) show similar observations for 100 and 500 random papers respectively. This observation motivates us to conduct detailed experiments by fixing papers in 10-year buckets and studying the citations they receive over time.

3.2 Fraction of citations in 10-year buckets

We now group all the papers in the Computer Science domain in different buckets as per publication years, with each bucket consisting of papers published in one decade. In Figure 2 (a), for each of these buckets, we plot the fraction of citations going to the current bucket and all the previous buckets. We note the following:

- The fraction of citations to the same bucket decreases over time (and those to all the older papers increase over time), consistent with the previous observations.
- If we consider the papers in a given bucket, the citations it receives decreases over the years. For instance, the papers in 1971-1980 received 70.5% of the citations in that decade but this

number reduces to 29.2, 6.4, 2.8 in the consecutive decades. This observation holds true for all the consecutive buckets as well.

To verify that these observations are universal, we take another dataset from the biomedical domain. We see similar observations with this dataset as well Figure 2(b).

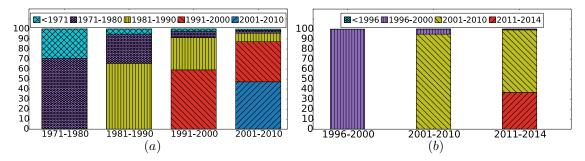


Figure 2: Stack plot showing citation distribution in 10-year buckets for papers in (a) Computer Science and (b) Biomedical domain, Each color/texture represents citations made to papers written in a time window, from future papers. Note that each such slab shrinks dramatically as time passes.

3.3 The case of different Computer Science fields

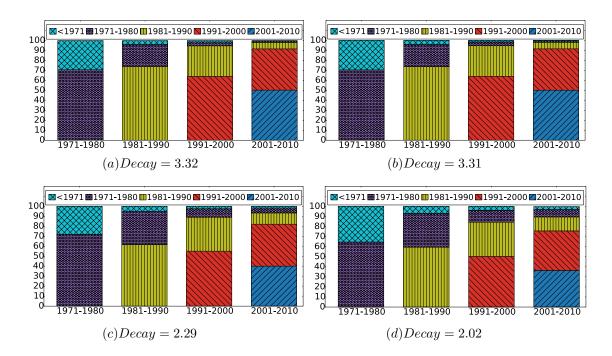


Figure 3: Stack plot showing citation distribution in 10-year buckets for (a) Data Mining, Machine Learning, Artificial Intelligence, Natural Language processing and Information Retrieval, (b) Distributed and Parallel Computing, Hardware and Architecture and real time and Embedded Systems, (c) Algorithms and Theory, Programming Languages and Software Engineering and (d) Algorithms,

We now repeat the same experiment by plotting the incoming citations to various sub-fields of computer science domain for various buckets. For instance, in Figure 3(a), the stacked plot at

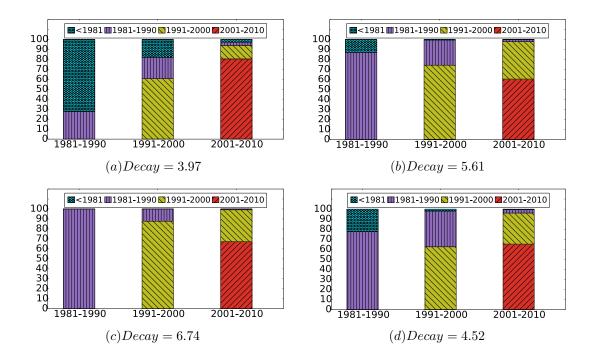


Figure 4: Stack plot showing citation distribution in 10-year buckets for (a) SIGMOD, (b) SIGIR, (c) ICML-NIPS and (d) VLDB-ICDE

1981-1990 denotes that among all the citations made to Data Mining etc. fields in that decade, what fractions of citations were made to the papers in the 1981-1990 (yellow), 1971-1980 (violet) and 1961-1970 (blue). We plot this for the following combinations of fields:

- Data Mining, Machine Learning, Artificial Intelligence, Natural Language processing and Information Retrieval Figure 3(a)
- Distributed and Parallel Computing, Hardware and Architecture and real time and Embedded Systems Figure 3(b)
- Algorithms and Theory, Programming Languages and Software Engineering Figure 3(c)

We see that earlier observations hold true for all these stacked bar charts as well, with fraction of citations going to older papers in that field decreasing over time. The amount by which this decay happens, though, is different for different fields (lower in case of Algorithms and Theory than others). If we study only the Algorithms field, the decay is even slower (Figure 3(d)).

To quantify this decay, we define a decay metric that measures the decrease in fraction of older citations for a field over the years. For papers published in a given time-interval, we compute the ratio of citations in consecutive year buckets and take a geometric mean over all such ratios. A high value of this geometric mean indicates a higher decay. So, in Figure 3, Algorithms and Theory field has the smallest decay. If we consider only the Algorithms field, decay is even smaller.

3.4 The case of different Computer Science Conferences

We further perform this analysis for various conferences in the field of Computer Science. Thus, for various time-points, we plot that among all the citations going to that particular conference, how many citations were made to the same time-point, previous one and so on. This analysis is performed for SIGMOD (Figure 4(a)), SIGIR (Figure 4(b)), ICML-NIPS (Figure 4(c)) and VLDB-ICDE (Figure 4(d)). The decay factors obtained for these conferences is much higher than observed for various fields in Figure 3. As a side remark, we observe that for SIGMOD, the fraction

of citations going to all the older papers decreases over time contrary to that reported in Verstak et al. [5]. However, our observation regarding the aging effect still holds true.

Further, we attempt to correlate the decay factor with the average value of 10 year impact factor [4] of the conferences in Table 2 and we find that they are negatively correlated, with high decay factor implying smaller Impact Factor.

Table 2: Correlation between decay and average value of 10-year impact factor.

Conf. Name	Decay Factor	Avg. IF $_{10}$
SIGMOD	3.97	3.50
SIGIR	5.61	2.77
ICML-NIPS	6.74	1.84
VLDB-ICDE	4.52	2.79

4 Discussions

A thorough analysis by fixing the set of cited papers in 10-year buckets helps us to reconcile the two contradictory views of the aggregate fate of ancient papers, that is, while the fraction of citations to all the papers older than a fixed number of years increases over time, that to a fixed set of old papers tends to decrease over time.

References

- [1] Tanmoy Chakraborty, Suhansanu Kumar, Pawan Goyal, Niloy Ganguly, and Animesh Mukherjee. Towards a stratified learning approach to predict future citation counts. In *Proceedings* of the 14th ACM/IEEE-CS Joint Conference on Digital Libraries, JCDL '14, pages 351–360. IEEE Press, 2014.
- [2] Tanmoy Chakraborty, Suhansanu Kumar, Pawan Goyal, Niloy Ganguly, and Animesh Mukherjee. On the categorization of scientific citation profiles in computer science. *Commun. ACM*, 58(9):82–90, August 2015.
- [3] P. Della Briotta Parolo, R. K. Pan, R. Ghosh, B. A. Huberman, K. Kaski, and S. Fortunato. Attention decay in science. *ArXiv e-prints*, March 2015.
- [4] Eugene Garfield. The history and meaning of the journal impact factor. *Jama*, 295(1):90–93, 2006.
- [5] Alex Verstak, Anurag Acharya, Helder Suzuki, Sean Henderson, Mikhail Iakhiaev, Cliff Chiung-Yu Lin, and Namit Shetty. On the shoulders of giants: The growing impact of older articles. CoRR, abs/1411.0275, 2014.