Original Articles

# Exploration and exploitation of Victorian science in Darwin's reading notebooks

Jaimie Murdock [a,b], Colin Allen [a,c,d], Simon DeDeo [a,b,e,f,*]

[a] Program in Cognitive Science, Indiana University, Bloomington, IN 47405, USA
[b] School of Informatics and Computing, Indiana University, 919 E. 10th Street, Bloomington, IN 47408, USA
[c] Department of History and Philosophy of Science and Medicine, Indiana University, Bloomington, IN 47405, USA
[d] School of Humanities and Social Sciences, Xi'an Jiaotong University, Xi'an, China
[e] Department of Social and Decision Sciences, Carnegie Mellon University, 5000 Forbes Avenue, BP 208, Pittsburgh, PA 15213, USA
[f] Santa Fe Institute, 1399 Hyde Park Road, Santa Fe, NM 87501, USA

A R T I C L E   I N F O

A B S T R A C T

Search in an environment with an uncertain distribution of resources involves a trade-off between exploitation of past discoveries and further exploration. This extends to information foraging, where a knowledge-seeker shifts between reading in depth and studying new domains. To study this decision-making process, we examine the reading choices made by one of the most celebrated scientists of the modern era: Charles Darwin. From the full-text of books listed in his chronologically-organized reading journals, we generate topic models to quantify his local (text-to-text) and global (text-to-past) reading decisions using Kullback-Liebler Divergence, a cognitively-validated, information-theoretic measure of relative surprise. Rather than a pattern of surprise-minimization, corresponding to a pure exploitation strategy, Darwin's behavior shifts from early exploitation to later exploration, seeking unusually high levels of cognitive surprise relative to previous eras. These shifts, detected by an unsupervised Bayesian model, correlate with major intellectual epochs of his career as identified both by qualitative scholarship and Darwin's own self-commentary. Our methods allow us to compare his consumption of texts with their publication order. We find Darwin's consumption more exploratory than the culture's production, suggesting that underneath gradual societal changes are the explorations of individual synthesis and discovery. Our quantitative methods advance the study of cognitive search through a framework for testing interactions between individual and collective behavior and between short- and long-term consumption choices. This novel application of topic modeling to characterize individual reading complements widespread studies of collective scientific behavior.

© 2016 Elsevier B.V. All rights reserved.

## 1. Introduction

The general problem of "information foraging" (Pirolli & Card, 1999) in an environment about which agents have incomplete information has been explored in many fields, including cognitive psychology (Hills, Todd, Lazer, Redish, & Couzin, 2015; Todd et al., 2012), neuroscience (Cohen, McClure, & Yu, 2007), economics (Azoulay-Schwartz, Kraus, & Wilkenfeld, 2004; March, 1991), finance (Uotila, Maula, Keil, & Zahra, 2009), ecology (Eliassen, Jørgensen, Mangel, & Giske, 2007; Stephens & Krebs, 1986), and computer science (Sutton & Barto, 1998). In all of these areas, the

searcher aims to enhance future performance by surveying enough of existing knowledge to orient themselves in the information space.

Individual scientists and scholars can be viewed as conducting a cognitive search (Todd et al., 2012) in which they must balance *exploration* of ideas that are novel to them against *exploitation* of knowledge in domains in which they are already expert (Berger-Tal, Nathan, Meron, & Saltz, 2014). Researchers have studied the exploration-exploitation trade-off in cognitive search at timescales of minutes up to years and decades. Laboratory experiments on visual attention are one example of this balancing act on short timescales (Chun & Wolfe, 1996), while studies of the recombination of patented technologies demonstrate long-term group behavior (Youn, Strumsky, Bettencourt, & Lobo, 2015). New advances in the digitization of historical archives enable longitudinal study of

how individuals explore and synthesize the work of their contemporaries and predecessors over the course of a lifetime.

As one of the most successful and celebrated scientists of the modern era, Charles Darwin's scientific creativity has been the subject of numerous narrative and qualitative studies (Gruber & Barrett, 1974; Johnson, 2010; Van Hulle, 2014). In part, these studies are possible because Darwin left his biographers careful records of his intellectual and personal life. These include records of the books he read from 1837 to 1860, a critical period which culminated in the publication of *The Origin of Species*. Table 1 summarizes key events in Darwin's life.

This article presents the first quantitative analysis of an important scientist's reading diaries, tracking how Darwin navigated the exploration-exploitation trade-off in choosing what to read. We link Darwin's reading records with the full text of the original volumes. We then use probabilistic topic models (Blei, Ng, & Jordan, 2003; Blei, 2012b) to represent the original text of each book Darwin read as a mixture of topics. We use information theory to measure the surprise, or unpredictability, of the next book that Darwin chose to read, compared to his past history of reading.

We present three key findings:

1. Darwin's reading patterns switch between both exploitation and exploration throughout his career. This is in contrast to a pure surprise-minimization strategy that consistently exploits content within a local region before moving on. The general trend, as Darwin's career develops, is towards increasing exploration.
2. In comparison to the publication order of the texts Darwin read, Darwin's reading order shows higher average surprise. This indicates that the order in which the books were written by the scientific community is less surprising than the order in which Darwin read them.
3. Darwin's strategies fall into three long-term epochs, or behavioral modes characterized by distinct patterns of surprise-seeking. These epochs correspond to three biographically significant periods: Darwin's post-*Beagle* studies, his extensive work on barnacles, and a final period leading to his synthesis of natural selection in the *Origin of Species*.

While the bulk of empirical studies in cognitive science are concerned with measuring population-level effects due to experimental manipulations, case studies play an important role in driving cognitive theorizing, experimentation, and modeling. For example, the case of the memory-impaired patient H.M. has driven many advances in cognitive neuroscience and computational models of memory (reviewed in Squire & Wixted (2011)). Other case studies, such as that of the frontal-lobe injury in Phineas Gage, provide important contrasts for later studies (reviewed in Macmillan (2000)).

Detailed longitudinal investigations of a single individual may involve repeated trials and changing strategies that cannot be observed in laboratory experiments involving a single task. Cognitive science could be enriched by using longitudinal studies such as ours to design laboratory studies with higher ecological validity. However, it is a challenging task to design laboratory studies of (for example) reading choices that take into account a subject's extensive prior history of reading decisions.

An advantage of taking Charles Darwin as a case study is the extensive attention that he has received from historians and biographers since his death. These studies provide a novel means for validating our mathematical tools. We will present a number of theoretical and empirical reasons why our methods measure cognitively-relevant features of a reader's experience. At the same time, the fact that our results also recover key features of Darwin's

**Table 1**

Timeline. Major events in Charles Darwin's life, including those marked in Fig. 1. This paper focuses on the critical period of his work from 1837 to 1860, leading to the publication of *The Origin of Species* (boundaries marked in bold). See Berra (2009) for an expanded chronology.

| *Major Events in Charles Darwin's Life (1809–1882)* | |
|---|---|
| 12 February 1809 | Born in Shrewsbury, England |
| 22 October 1825 | Matriculates at University of Edinburgh |
| 15 October 1827 | Admitted to Christ's College, Cambridge |
| 27 December 1831 | Departs England aboard the *HMS Beagle* |
| 2 October 1836 | Return to England aboard the *HMS Beagle* |
| **July 1837** | **First entries in reading notebooks** |
| August 1839 | Publication of *The Voyage of the Beagle* (1st edition) |
| May 1842 | Writes the 1st Essay on Species |
| 4 July 1844 | Writes the 2nd Essay on Species |
| August 1845 | Publication of *The Voyage of the Beagle* (2nd edition) |
| 1 October 1846 | Begins barnacle project |
| 19 February 1851 | Publishes first volume of barnacle work |
| 9 September 1854 | Begins sorting notes on natural selection |
| 14 May 1856 | Starts writing "large work" on species |
| 24 November 1859 | Publication of *The Origin of Species* (1st edition) |
| **13 May 1860** | **Last entry in reading notebooks** |
| 24 February 1871 | Publication of *The Descent of Man* |
| 19 February 1872 | Publication of *The Origin of Species* (6th and final edition) |
| 21 April 1882 | Dies at Down House in Kent, England |

intellectual life provides additional support for the reliability of our methods.

We cannot claim Darwin's information foraging behavior is typical for either scientists or the population-at-large; however, our analytic methods may be applied to other case studies to find population-level generalizations, provided there is access to the full text of each reading.

Our approach contrasts with previous uses of topic modeling to analyze the large-scale structure of scientific disciplines (Blei & Lafferty, 2007; Cohen Priva & Austerweil, 2015; Griffiths & Steyvers, 2004; Hall, Jurafsky, & Manning, 2008) and the humanities (Blei, 2012a; Jockers, 2013; Mohr & Bogdanov, 2013), which are each created through the collective effects of individual-level behavior. Previous models of historical records have focused on language use as an indication of larger shifts in style (Hughes, Foti, Krakauer, & Rockmore, 2012; Underwood & Sellers, 2012), learnability (Hills & Adelman, 2015), or content (Goldstone & Underwood, 2014; Klingenstein, Hitchcock, & DeDeo, 2014; Michel et al., 2011) of significant portions of publications in a field, including a study of *Cognition* itself (Cohen Priva & Austerweil, 2015).

These works model the collective state of all published works at a particular date, but obscure the role of individual foraging behavior. By focusing on a single individual for whom ample records exist, we gain access to what Tria, Loreto, Servedio, and Strogatz (2014) describe as "the interplay between individual and collective phenomena where innovation takes place".

## 2. Materials and methods

### 2.1. Darwin's reading notebooks

Darwin was a meticulous record-keeper—starting in April 1838, he kept a notebook of "books to be read" and "books read". These records span the 23 years from 1837 to 1860, tracking his reading choices from just after his return to England aboard the *HMS Beagle* to just after the publication of *The Origin of Species*. We located the full-text of 665 of the 687 (96.7%) English non-fiction mentioned in these reading notebooks through a variety of online digital libraries. See Appendix A for details of corpus curation.