

Northeastern University





CESS $() + R() \times A$

Xindi Wang Northeastern University, Boston, USA | Center for Complex Network Research (CCNR)

Northeastern University

Network Science Institute

Barabási Lab





More than 3 Million books are published every year in the US.



More than 3 Million books are published every year in the US.

GIT II

Only 500 of those 3 Million become bestsellers!



Why some books become successful while others fail?



ALL THE DATA!!!







- Covers approximately 85% of print trade book sales, over 500,000 ISBNs tracked in a week for 15 years
- Author's career history







- Covers approximately 85% of print trade book sales, over 500,000 ISBNs tracked in a week for 15 years
- Author's career history





- Best seller lists with many categories
- Updated weekly from 2008
- Well known measure of success for books



- Covers approximately 85% of print trade book sales, over 500,000 ISBNs tracked in a week for 15 years
- Author's career history



- Detailed info of books and authors
- Ratings and reviews of books and authors
- Readers data: books read, friend, etc.



- Best seller lists with many categories
- Updated weekly from 2008
- Well known measure of success for books

npd Bookscan The New York Times



- Covers approximately 85% of print trade book sales, over 500,000 ISBNs tracked in a week for 15 years
- Author's career history



- Detailed info of books and authors
- Ratings and reviews of books and authors
- Readers data: books read, friend, etc.



- Best seller lists with many categories
- Updated weekly from 2008
- Well known measure of success for books

The New York Times





- Page-views: how many clicks for a page
- Indicates global interest in an author/book
- Info about famous books & authors





• Bestsellers study, specifically sales pattern

Predict success before the book published

OUTLINE



SALES PATTERN





Most books follow a universal "quick rise, slower decay" pattern Peak arrives within 10 weeks after publishing

SALES PATTERN





SALES PATTERN MODELING

- •
- likely it is to be purchased again.
- Aging: $A_i(t) = \frac{1}{\sqrt{2\pi\sigma_i t}} \exp\left[-\frac{(\ln t \mu_i)^2}{2\sigma_i^2}\right]$ σ_i : the decay rate capturing the longevity.

Dashun Wang, Chaoming Song, and Albert-László Barabási, "Quantifying long-term scientific impact," Science 342, 127–132 (2013).

Fitness η_i : the book's ability to respond to the taste of a wide readership.

Preferential attachment: the higher the up-to-date sales of a book S_i^t , the more

 μ_i : the book's immediacy determined by the time the sales reach its peak



 λ_i is the relative fitness proportional to η_i







PREDICT THE FUTURE $S_i^{\infty} = m(e^{\lambda_i} - 1)$











• Sales pattern follows a "quick rise, slow decay" pattern, the peak usually arrives within 10 weeks.

- usually arrives within 10 weeks.
- lot of sales data

• Sales pattern follows a "quick rise, slow decay" pattern, the peak

• We can infer future sales from past sales, but this inference needs a



- usually arrives within 10 weeks.
- lot of sales data

More investigations about bestseller books and authors'

Success in Books: A Big Data Approach to Bestsellers, Burcu Yucesoy, Xindi Wang, Junming Huang, Albert-László Barabási, EPJ Data Science (Submitted)

• Sales pattern follows a "quick rise, slow decay" pattern, the peak

• We can infer future sales from past sales, but this inference needs a



Can we predict sales before the book is published?





- Fame
- Publishing History



- Fame
- Publishing History



Famous because he was president





- Fame
- Publishing History



Famous because he was president Famous because of her books





- Fame
- Publishing History





- Fame
- Publishing History



FEATURES

Book

Publishing Month





- Fame
- Publishing History



FEATURES



Book

Publishing Month

Publisher IIII

• Publisher (Imprint) Prominence



WIKIPEDIA The Free Encyclopedia

• Fame

Author

• Publishing History

• Fame

- Cumulative Fame
- Fame Establishing Days
 Previous sales in this genre
- Fame
- Recent Fame

FEATURES





Total previous sales

Publishing History

- Normalized Cumulative
 Previous sales in other genres
 - Normalized previous sales









- Topic •
- **Publishing Month** •







- Genre
- Topic
- Publishing Month



- Genre
 - Bisac Code





FIC031000 BIO003000





- Genre
- Topic
- **Publishing Month**



- Genre
 - Bisac Code





FIC031000 BIO003000







FEATURES

Nonnegative matrix factorization (NMF) on book summary







- Genre
- Topic
- **Publishing Month**



- Genre
 - Bisac Code





FIC031000 BIO003000







FEATURES

Nonnegative matrix factorization (NMF) on book summary





Publishing Month







• Publisher (Imprint) Prominence







- Linear regression?
 - Problem of linear regression: unbalanced data

HOW TO BUILD THE MODEL?



- Linear regression?
 - Problem of linear regression: unbalanced data



HOW TO BUILD THE MODEL?

Systematic Underprediction



their sales, where would we place a new book in this sequence?

Book 2 Sale

Book 1 Sale

Given a sequence of previously published books ranked by





LEARNING TO PLACE

Pairwise Comparison

Placing objects



TRAINING PHASE









There maybe conflicts among predicted labels





































Nonfiction

RESULTS







Predicted Sale is the middle point of most voted interval



Fiction



RESULTS

Nonfiction





CONCLUSION



CONCLUSION

• We build a model Learning to Place that could predict the success of books pretty accurately



CONCLUSION

- We build a model Learning to F books pretty accurately
- We know the feature importar and Nonfiction)

• We build a model Learning to Place that could predict the success of

• We know the feature importance for different categories (Fiction



WHAT'S NEXT



Improvement on Nonfiction

WHAT'S NEXT



- Improvement on Nonfiction

WHAT'S NEXT

• Does social media play a role in publishing: Information cascading



- Improvement on Nonfiction
- Author career

WHAT'S NEXT

• Does social media play a role in publishing: Information cascading





Northeastern University Network Science Institute



Burcu Yucesoy



Onur Varol



Tina Eliassi-Rad

Barabási Lab







Junming Huang



Albert-László Barabási





ESTIONS?

Email: xindi.w1993@gmail.com



SALES PATTERN MODELING

- **Fitness** η_i : the book's ability to respond to the taste of a wide readership. •
- likely it is to be purchased again.
- Aging: $A_i(t) = \frac{1}{\sqrt{2\pi\sigma_i t}} \exp\left[-\frac{(\ln t \mu_i)^2}{2\sigma_i^2}\right]$ μ_i : the book's immediacy determined by the time the sales reach its peak
 - σ_i : the decay rate capturing the longevity.

The probability of a book i to be purchased at a time t after publication $\Pi_i(t) \sim \eta_i S_i^t A_i$

Dashun Wang, Chaoming Song, and Albert-László Barabási, "Quantifying long-term scientific impact," Science 342, 127–132 (2013).



Preferential attachment: the higher the up-to-date sales of a book S_i^t , the more

HOW TO CALCULATE ROC



At each threshold, how many false positive and true positive