

Temporal patterns in information and social systems

Matúš Medo

University of Fribourg, Switzerland

COST Workshop

Quantifying scientific impact: networks, measures, insights?

12-13 February, 2015, Zurich

Outline

- 1 Growing networks with fitness and aging
- 2 Temporal bias of PageRank
- 3 Leaders and followers and the consequences

The common theme

Temporal patterns and the role of time in information and social systems.

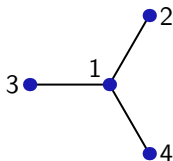
Preferential attachment (PA)

- A classical network model
 - Yule (1925), Simon (1955), Price (1976), Barabási & Albert (1999)
 - Growth of cities, citations of scientific papers, WWW,...

Preferential attachment (PA)

- A classical network model
 - Yule (1925), Simon (1955), Price (1976), Barabási & Albert (1999)
 - Growth of cities, citations of scientific papers, WWW,...
- Nodes and links are added with time
- Probability that a node acquires a new link proportional to its current degree

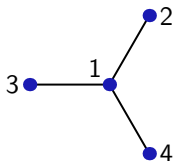
$$P(i, t) \sim k_i(t)$$



Preferential attachment (PA)

- A classical network model
 - Yule (1925), Simon (1955), Price (1976), Barabási & Albert (1999)
 - Growth of cities, citations of scientific papers, WWW,...
- Nodes and links are added with time
- Probability that a node acquires a new link proportional to its current degree

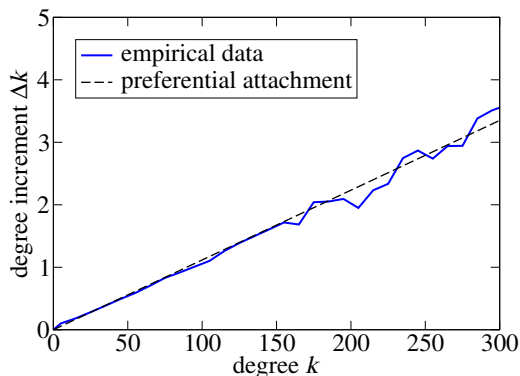
$$P(i, t) \sim k_i(t)$$



- Pros: simple, produces a power-law degree distribution

PA in scientific citation data

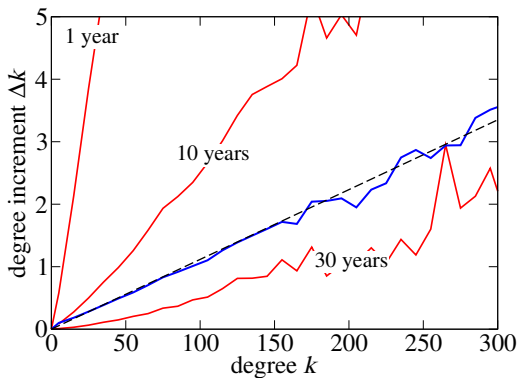
Journals of the American Physical Society from 1893 to 2009:



See also Adamic & Huberman (2000), Redner (2005), Newman (2009),...

PA in scientific citation data

Journals of the American Physical Society from 1893 to 2009:



Time decay is fundamental

"All the News That's Fit to Print."

The New York Times.

THE WEATHER.

THE. 5.61. NO. 23.000. NEW YORK, FRIDAY, APRIL 15, 1912. TWENTY-FIVE CENTS. ONE CENT. LONDON 2P. 11. 10. 11. 12. 13. 14. 15. 16. 17. 18. 19. 20. 21. 22. 23. 24. 25. 26. 27. 28. 29. 30. 31. 32. 33. 34. 35. 36. 37. 38. 39. 40. 41. 42. 43. 44. 45. 46. 47. 48. 49. 50. 51. 52. 53. 54. 55. 56. 57. 58. 59. 60. 61. 62. 63. 64. 65. 66. 67. 68. 69. 70. 71. 72. 73. 74. 75. 76. 77. 78. 79. 80. 81. 82. 83. 84. 85. 86. 87. 88. 89. 90. 91. 92. 93. 94. 95. 96. 97. 98. 99. 100.

TITANIC SINKS FOUR HOURS AFTER HITTING ICEBERG; 866 RESCUED BY CARPATHIA, PROBABLY 1250 PERISH; ISMAY SAFE, MRS. ASTOR MAYBE, NOTED NAMES MISSING

Col. Astor and Bride, Isaac Straus and Wife, and Maj. Butt Aboard.

"HOLE OF DEEP" FOLLOWED

Women and Children Put Safe in Lifeboats and Are Expected to Be Safe on Carpathia.

PICKED UP AFTER 8 HOURS

Investigator Called White Star Ship for News of His Father and Lovers Missing.

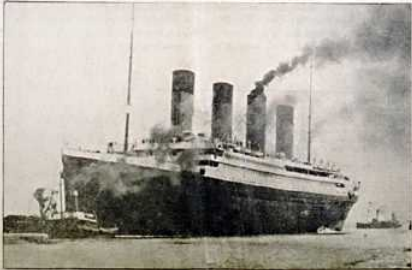
FRANKLIN HOSPITAL ALL SET

Manager of the Life Rafted Thence After Stranahan Goes After Him and One Dies.

HEAD OF THE LINE ANCHORED

A Boat Seen Heading West This Morning With Two Men to Rescue at Once.

The Atlantic, that the Titanic, the largest ship in the world, was sunk at an estimated rate of 100 miles per hour, and that the vessel of the American coastwise line, the Carpathia, was the first to reach the scene of the disaster.



Biggest Liner Plunges to the Bottom at 2:20 A. M.

RESCUERS THERE TOO LATE

Except to Pick Up the Few Survivors Who Took to the Lifeboats.

WOMEN AND CHILDREN FIRST

General Serpa's Boats to Save Them with the Survivors.

SEEK SEARCH FOR OTHERS

The Carpathia Starts by an Effort of Picking Up Other Boats or Rafts.

CLIPPING SENDS THE NEWS

City Will to Hear 'TITANIC' Message to Be Sent to the Ship.

LARGE REPORTS MADE

REUTERS, APRIL 15, 1912. NEW YORK. (APRIL 15, 1912.)—The Titanic, the largest ship in the world, was sunk at an estimated rate of 100 miles per hour, and that the vessel of the American coastwise line, the Carpathia, was the first to reach the scene of the disaster.

Growing networks with fitness and aging

(PRL 107, 238701, 2011)

- Probability that node i attracts a new link

$$P(i, t) \sim \underbrace{k_i(t)}_{\text{degree}} \times \underbrace{f_i}_{\text{fitness}} \times \underbrace{D_R(t)}_{\text{aging}}$$

relevance

- The aging factor $D_R(t)$ decays with time: a decay of relevance
- When $D_R(t) \rightarrow 0$, the popularity of nodes eventually saturates

Growing networks with fitness and aging

(PRL 107, 238701, 2011)

- Probability that node i attracts a new link

$$P(i, t) \sim \underbrace{k_i(t)}_{\text{degree}} \times \underbrace{f_i}_{\text{fitness}} \times \underbrace{D_R(t)}_{\text{aging}}$$

relevance

- The aging factor $D_R(t)$ decays with time: a decay of relevance
- When $D_R(t) \rightarrow 0$, the popularity of nodes eventually saturates
- The bottom line:
 - **Good:** Produces various realistic degree distributions (power-law, etc.)
 - **Bad:** Difficult to validate (high-dimensional statistics)
 - **Good:** This model explains the data much better than any other (PRE 89, 032801, 2014)

Growing networks with fitness and aging

(PRL 107, 238701, 2011)

- Probability that node i attracts a new link

$$P(i, t) \sim \underbrace{k_i(t)}_{\text{degree}} \times \underbrace{f_i}_{\text{fitness}} \times \underbrace{D_R(t)}_{\text{aging}}$$

relevance

- The aging factor $D_R(t)$ decays with time: a decay of relevance
- When $D_R(t) \rightarrow 0$, the popularity of nodes eventually saturates
- Point to note:

Paper popularity grows exponentially with fitness (quality)



Fitness depends logarithmically on popularity

Two forms of aging in information networks

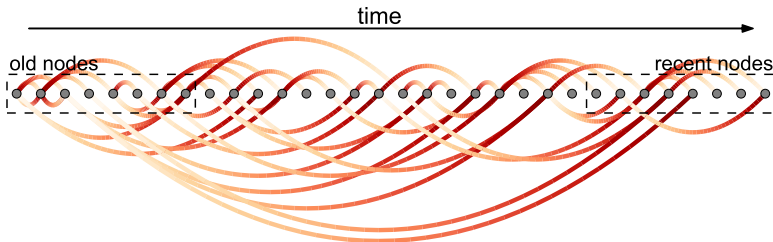
- The decay of relevance: $D_R(t)$
 - Node relevance influences the in-coming links

Two forms of aging in information networks

- The decay of relevance: $D_R(t)$
 - Node relevance influences the in-coming links
- The decay of activity: $D_A(t)$
 - Nodes activity influences the out-going links
 - Activity decays in time (mostly)

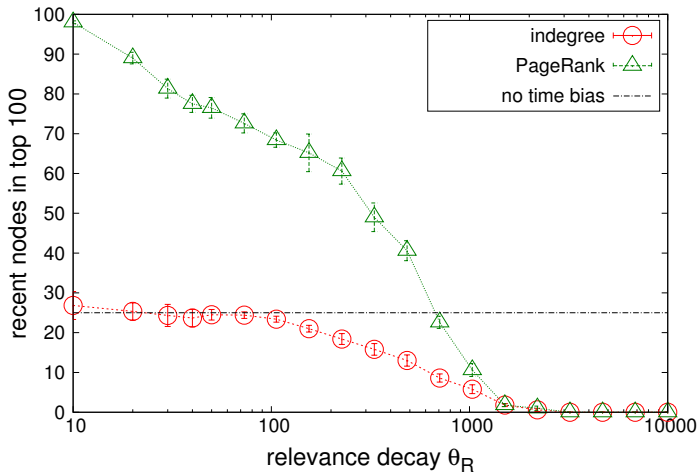
Two forms of aging in information networks

- The decay of relevance: $D_R(t)$
 - Node relevance influences the in-coming links
- The decay of activity: $D_A(t)$
 - Nodes activity influences the out-going links
 - Activity decays in time (mostly)

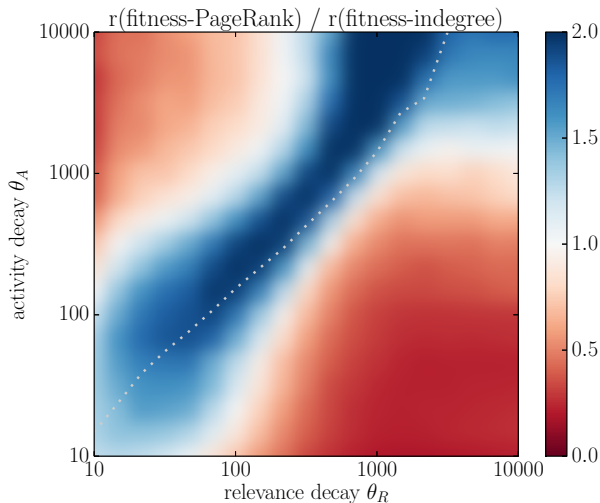


A growing network with a quick decay of attractiveness and no decay of activity

The biases of PageRank



The biases of PageRank



The biases of PageRank

Implications

In citation data, the time scales of relevance and activity decay are very different ($\Theta_A = 0$ because outgoing links are created only upon arrival).

PageRank (and its variants) is nevertheless commonly applied on citation data. One should think twice!

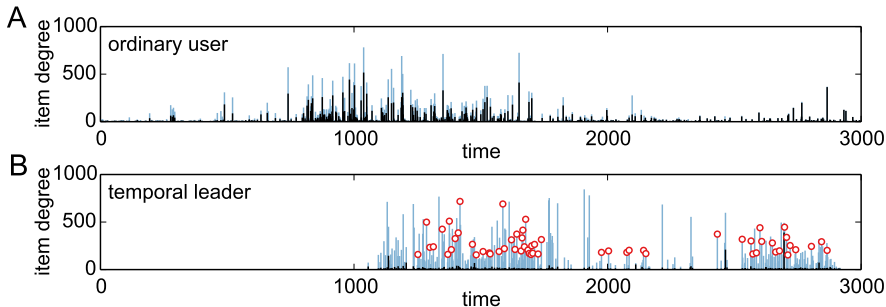
The case for leaders in social systems

- Bipartite user-item data (*who* bought *what* at Amazon.com)
 - Similar behavior in monopartite social data (user-user)
- Most users are driven by item popularity (*followers*)
- Some users are driven by item fitness (*leaders*)

The case for leaders in social systems

- Bipartite user-item data (*who bought what* at Amazon.com)
 - Similar behavior in monopartite social data (user-user)
- Most users are driven by item popularity (*followers*)
- Some users are driven by item fitness (*leaders*)
- A user makes a *discovery* when they are among the first 5 users to collect an eventually highly popular item (top 1% of all items are used as target)
- A new metric, *user surprisal*, shows that there are users who make discoveries so often that it cannot be explained by luck

Leaders in Amazon data



Black bars: popularity of collected items when they are collected.

Blue bars: final popularity of collected items.

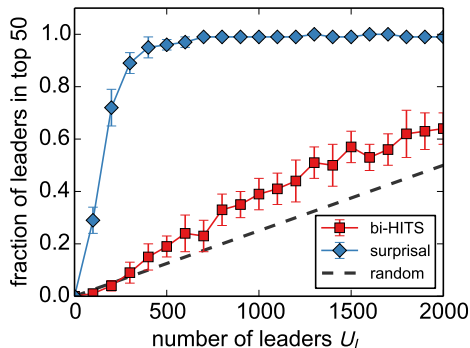
Red circles: discoveries.

The game changer

- Network growth model with two rules reproduces the real data patterns
 - 1 Followers choose items driven by $k_i(t)D_R(t)$
 - 2 Leaders choose items driven by $f_i(t)D_R(t)$

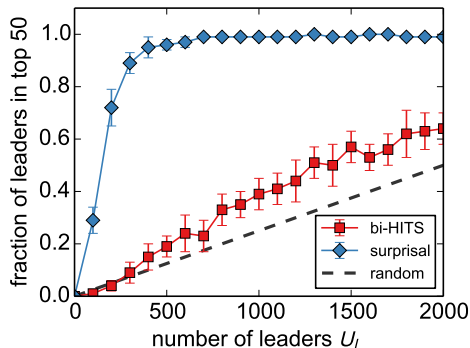
The game changer

- Network growth model with two rules reproduces the real data patterns
 - 1 Followers choose items driven by $k_i(t)D_R(t)$
 - 2 Leaders choose items driven by $f_i(t)D_R(t)$
- Model data poses a puzzle to classical ranking algorithms



The game changer

- Network growth model with two rules reproduces the real data patterns
 - 1 Followers choose items driven by $k_i(t)D_R(t)$
 - 2 Leaders choose items driven by $f_i(t)D_R(t)$
- Model data poses a puzzle to classical ranking algorithms



Reason: Insightful choices of the leaders are copied by the followers. All users ultimately collect items of the same fitness and an algorithm acting on a static data snapshot cannot distinguish them.

Solution: Algorithms that take time into account adequately.

Coming back to the guiding questions:

“Are the implicit assumptions of centrality measures justified in scientometrics?”

“It doesn’t seem to be the case!”

“Are altmetrics shallow?”

“Build on the community structure of science!”

(PLoS ONE 9, e112022, 2014)

Thank you for your attention