

Web

Master 1 IFI



Andrea G. B. Tettamanzi
Université Nice Sophia Antipolis
Département Informatique
andrea.tettamanzi@unice.fr

Unit 8

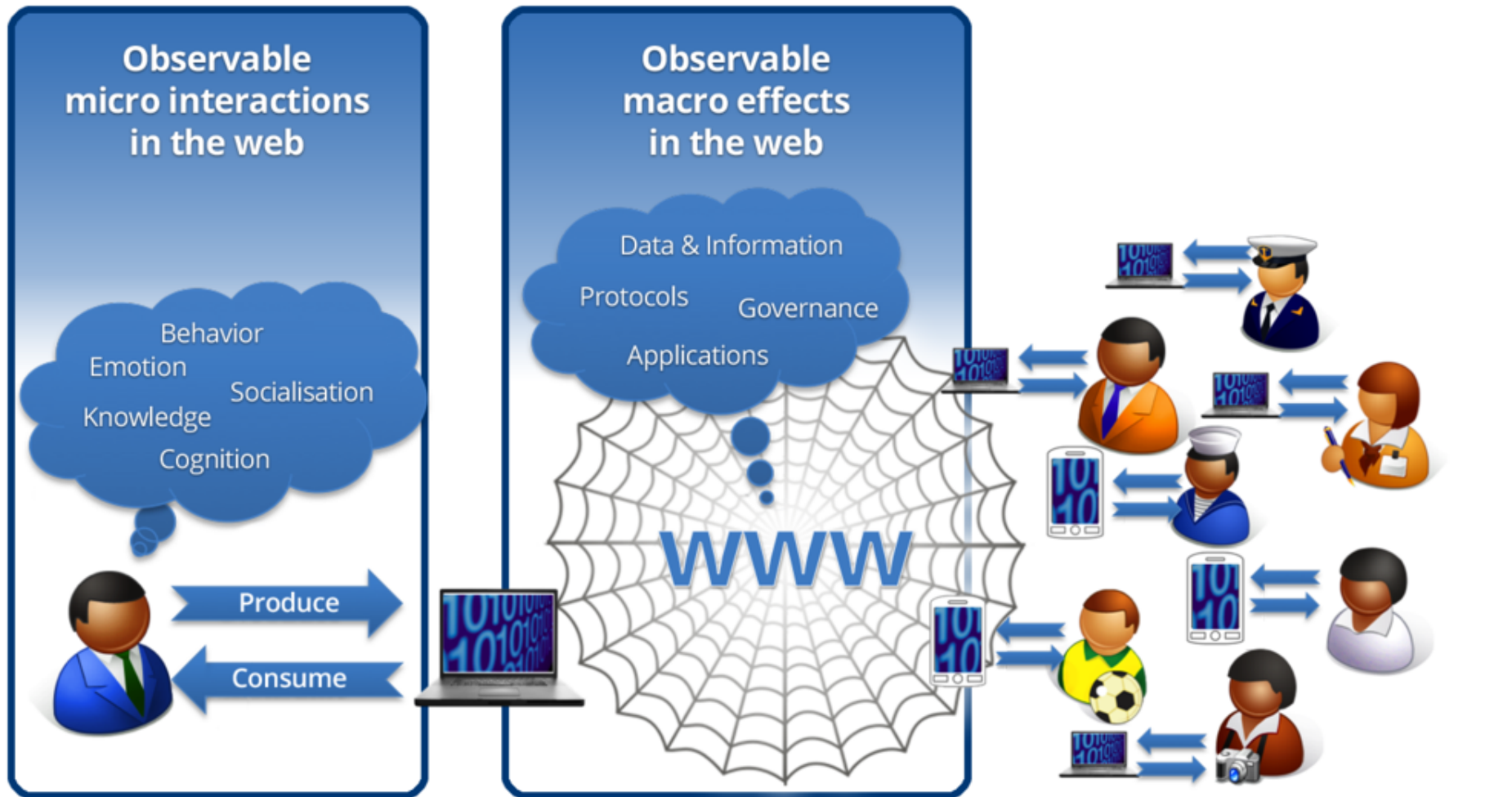
Web Science

–

PageRank and how Google turns words into money

Web Science

- A deliberately ambiguous phrase.
 - Physical Science is an analytic discipline that aims to find laws that generate or explain observed phenomena
 - Computer Science is predominantly (though not exclusively) synthetic, in that formalisms and algorithms are created in order to support particular desired behavior
- Web science is a merging of these two paradigms
 - The Web needs to be studied and understood
 - ...and it needs to be engineered.
- At the micro scale, the Web is a piece of engineering
- At the macro scale, the Web is an emergent phenomenon
- Interdisciplinary: CS + Mathematics + Sociology + Economics



Web Science Topics

- Modeling Web-related structures, data, users and behaviors
- Analysis of online social and information networks, social media analysis
- Social machines, crowd computing, collective intelligence, and collaborative production (e.g., prediction markets)
- Web Economics
- Sentiment Analysis and Opinion Mining
- Legal Aspects of the Web
- Ethical Challenges
- Web access, literacy, divides, inclusion, exclusions, and development
- Humanities, arts, and culture on the Web

Introduction to PageRank

- Key statistics about Alphabet Inc. (= Google), as of April 9, 2019
 - Market capitalization: ~ \$838 billion (2014: \$375 billion)
 - Revenue : \$136 billion (2014: \$62 billion)
 - EBITDA : \$34.9 billion (2014: \$18.6 billion) \$1,107/s !!!
 - Full-time employees: 98,771 (2014: 54,000)
- As a comparison:
 - GDP of Ukraine: ~ \$126.4 billion
 - If Google were a country, it would be 57th by GDP out of 186
 - In 2016, Alphabet was 94th among the world's corporations by capitalization and 2nd among publicly traded companies
- Not bad for “just” a search engine...

The Key of Success

- Google's success is based on two algorithms :
 - **PageRank**
 - **AdWords + AdSense**
- The former allows Google to rank search results:
 - It gives Google its **use value**
 - It has imposed Google as a market leader
- The latter generates the impression of advertisements targeted on the interests of the audience of a Web page:
 - It gives Google its **exchange value**
 - AdWords allows buying traffic, AdSense allows selling traffic

Agenda

- PageRank
- AdWords + AdSense
- Lab work

Part I

PageRank

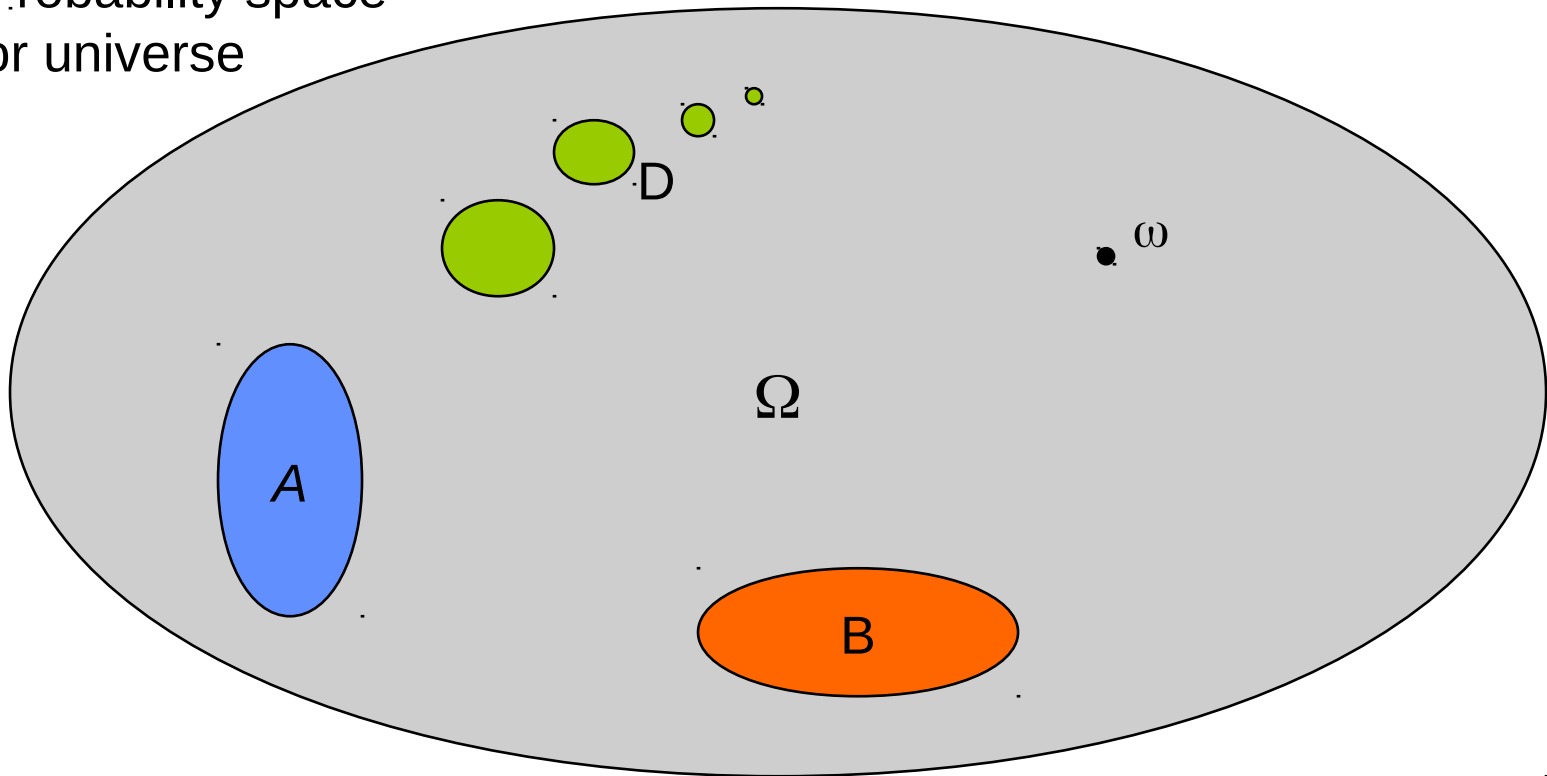


Basic Intuition

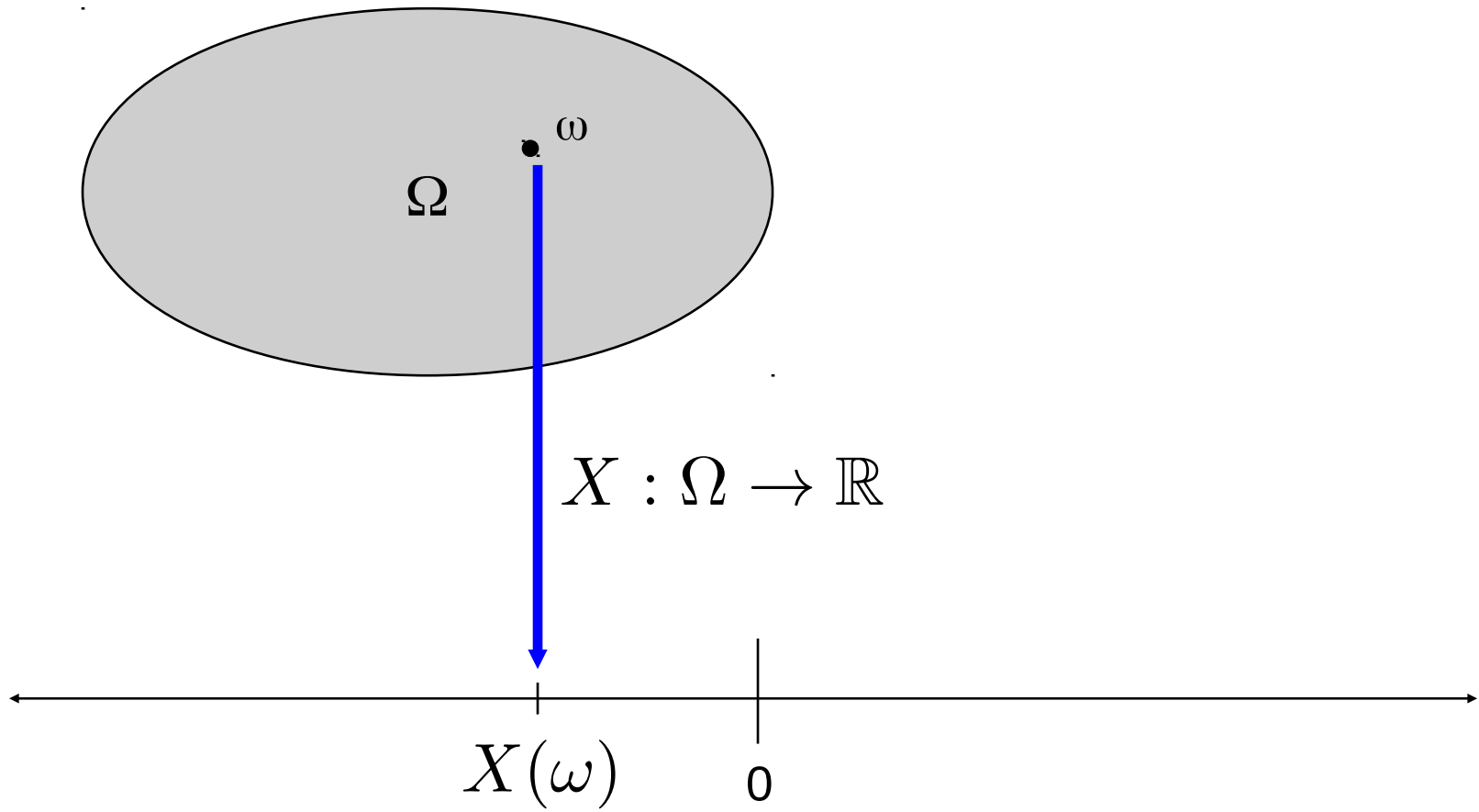
- The WWW as a directed graph
 - Its nodes are the HTML pages
 - Its arcs are the ` . . . ` hyperlinks
- Which pages would a **random surfer** visit?
 - The random surfers would start at a random page
 - They would jump from one page to the next by clicking a random hyperlink
 - Idea: measure the importance of a page by the probability that it is visited at time t by a random surfer!
- This probability is the visit frequency of the page

Events

Probability space
or universe



Random Variables



Random Processes

A sequence of random variables

$$X_1, X_2, \dots, X_t, \dots$$

Each equipped with its own probability distribution.

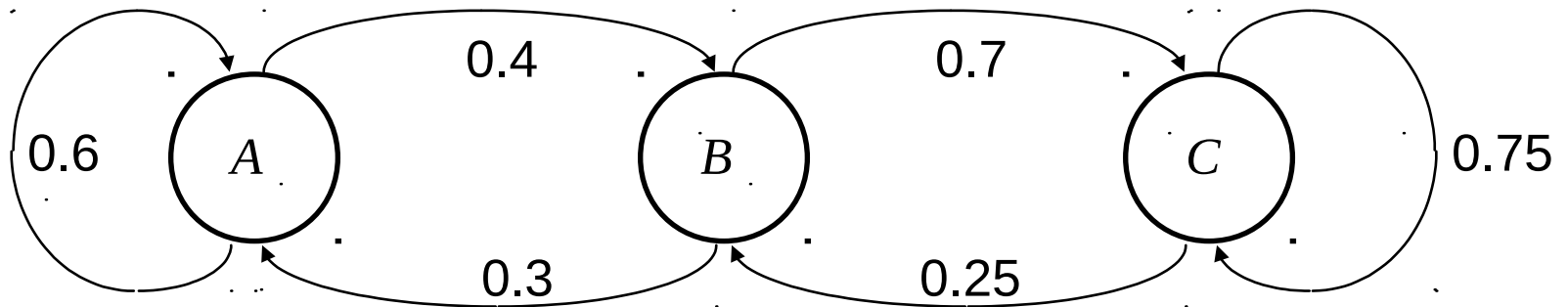
Notation: $\{X_t(\omega)\}_{t=0,1,\dots}$

Markov Chains

A random process $\{X_t(\omega)\}_{t=0,1,\dots}$

is a Markov chain if and only if, for all t ,

$$\Pr[X_t = x \mid X_0, X_1, \dots, X_{t-1}] = \Pr[X_t = x \mid X_{t-1}]$$



Transition Matrix

$$\mathbf{T} = \begin{bmatrix} \Pr(X_t = x_1 \mid X_{t-1} = x_1) & \dots & \Pr(X_t = x_n \mid X_{t-1} = x_1) \\ \Pr(X_t = x_1 \mid X_{t-1} = x_2) & \dots & \Pr(X_t = x_n \mid X_{t-1} = x_2) \\ \vdots & & \vdots \\ \Pr(X_t = x_1 \mid X_{t-1} = x_n) & \dots & \Pr(X_t = x_n \mid X_{t-1} = x_n) \end{bmatrix}$$

\mathbf{T} is a stochastic matrix:

$$\forall i, \quad \sum_{j=1}^n \Pr(X_t = x_j \mid X_{t-1} = x_i) = 1$$

“Idealized” Definition of PageRank

$q_i = \#$ outgoing links from page i

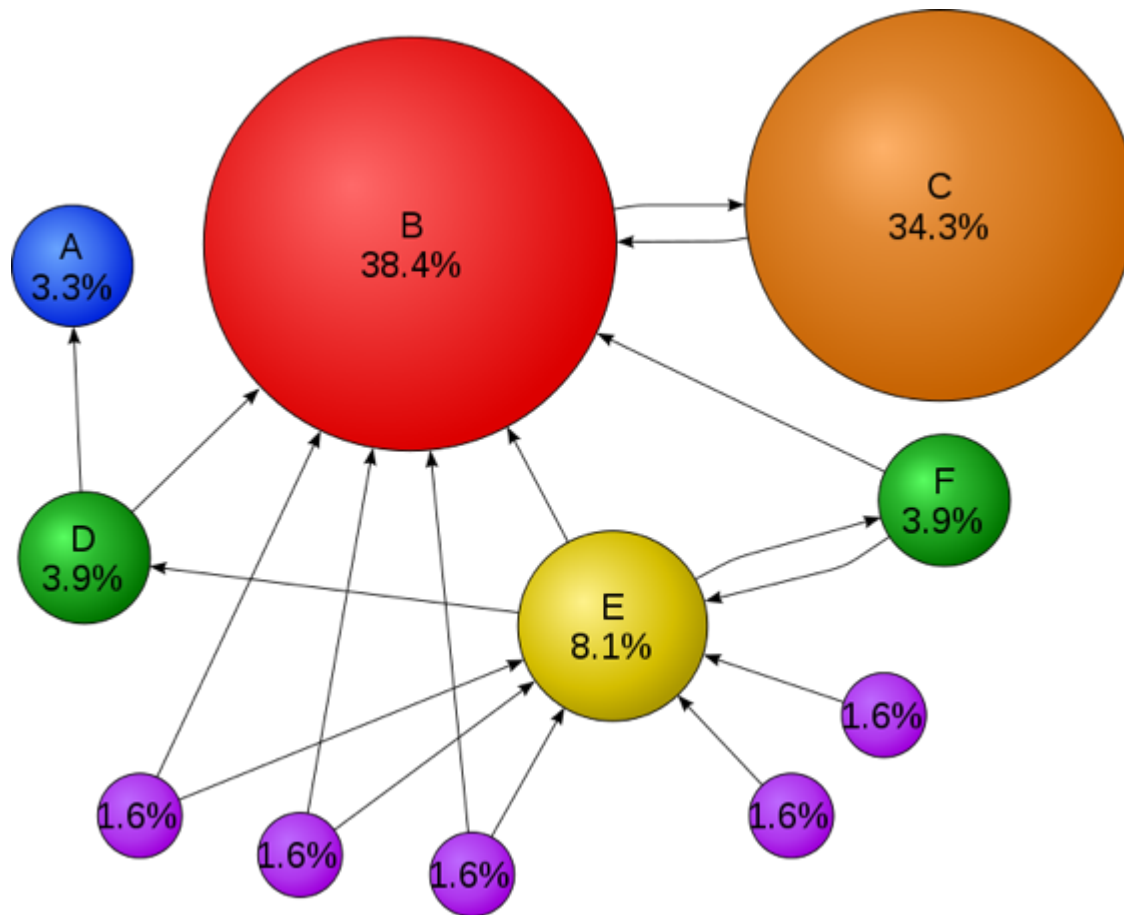
$$\mathbf{H} = (h_{ij})$$

$$h_{ij} = \begin{cases} 1/q_i & \text{there exists a link from } i \text{ to } j; \\ 0 & \text{otherwise.} \end{cases}$$

The diagram consists of two rectangular boxes connected by a large blue double-headed arrow. The left box contains the equation $\pi_j = \sum_i \pi_i h_{ij}$. The right box contains the matrix equation $\pi = \pi \mathbf{H}$.

$$\pi_j = \sum_i \pi_i h_{ij} \quad \longleftrightarrow \quad \pi = \pi \mathbf{H}$$

Example



Basic Hypothesis

A Web page is important
insofar as it is referenced by
other important pages

Analysis of the Definition

- There are three factors that determine the PageRank of a page:
 - The number of links pointing towards it;
 - The propensity of the pages containing those links to direct surfers towards it, i.e., the total number of outgoing links;
 - The PageRank of the pages containing those links
- The idealized model has two problems:
 - Pages without outgoing links (*dangling pages*), which can capture surfers.
 - A surfer may also get trapped in a *bucket*, a reachable and strongly connected component, without outgoing arcs towards the rest of the graph.

Real Model: the Google Matrix

- The lines of matrix \mathbf{H} having all zero elements, corresponding to pages without outgoing links, are replaced by a uniform or arbitrary distribution.
- Let \mathbf{S} be the matrix thus modified.
- To solve the problem with *buckets*, Brin and Page propose to replace matrix \mathbf{S} by the Google matrix:

$$\mathbf{G} = \delta \mathbf{S} + (1 - \delta) \mathbf{E}$$

damping factor δ \leftarrow *Teleportation matrix* \mathbf{E}

$$\mathbf{E} = \begin{bmatrix} 1/n & 1/n & \cdots & 1/n \\ \vdots & \vdots & & \vdots \\ 1/n & 1/n & \cdots & 1/n \end{bmatrix}$$

Interpreting the Google Matrix

- The definition of the Google matrix may be explained as follows
 - With probability δ , the random surfer follows the next link
 - With probability $1 - \delta$, the random surfer gets tired following links and directs the browser to a novel URL, which has nothing to do with the current page.
 - In this case, the surfer is “teleported” to this novel page
- The inventors of PageRank suggest a damping factor $\delta = 0.85$:
 - On average, after following 5 links, the surfer chooses a new random page.
- The PageRank vector is therefore π such that

$$\pi = \pi \mathbf{G}$$

Existence and Uniqueness of the PageRank vector

- The π vector is an eigenvector of \mathbf{G} of eigenvalue 1.
- The \mathbf{S} matrix is stochastic, as is matrix \mathbf{E} .
- The \mathbf{G} matrix is, therefore, stochastic as well.
- If \mathbf{G} is stochastic, equation $\pi = \pi\mathbf{G}$ has at least one solution.
- According to Perron-Frobenius' Theorem, if \mathbf{A} is an irreducible non-negative square matrix, then there exists a vector \mathbf{x} such that $\mathbf{x}\mathbf{A} = r\mathbf{x}$, where r is the spectral radius of \mathbf{A} .
- The \mathbf{S} matrix is likely to be reducible; however, thanks to the teleportation matrix, \mathbf{G} is certainly irreducible.
- Furthermore, since \mathbf{G} is stochastic, its spectral radius is 1.
- As a consequence, a PageRank vector > 0 exists and is unique.

PageRank and Markov Theory

- The random walk model on the Web graph, modified with teleportation, naturally induces a Markov chain with a finite (albeit huge) number n of states (= pages)
- \mathbf{G} is the transition matrix of such Markov chain
- Since \mathbf{G} is irreducible, the chain is ergodic and it has a unique stationary distribution, corresponding to the PageRank vector π .

Computing the PageRank Vector (1)

- The **power method** is a numerical method which allows to determine the greatest (in absolute value) eigenvalue of a matrix with real coefficients.
- We take a random vector \mathbf{x} and we compute the recurrence:

$$\mathbf{x}^{(0)} = \mathbf{x}, \quad \mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} \mathbf{A} / \|\mathbf{A}\|$$

- This sequence converges to the greatest (in absolute value) eigenvalue of matrix \mathbf{A}
- To compute π , we start from vector $\mathbf{u} = (1/n, \dots, 1/n)$ and we stop as soon as

$$\|\pi^{(t+1)} - \pi^{(t)}\| < \epsilon$$

Computing the PageRank Vector (2)

- The convergence speed of the power method applied to matrix \mathbf{G} is of the same order as the rate by which δ^k goes to 0.
- For instance, for $\delta = 0.85$:
 - 43 iterations \rightarrow precision of 3 decimal digits
 - 142 iteration \rightarrow precision of 10 decimal digits
- We also observe that the power method applied to matrix \mathbf{G} can be expressed in terms of matrix \mathbf{H}
- \mathbf{H} is an extremely sparse matrix, which can be stored in a memory space of size $O(n)$
- According to rumors, Google recomputes π once per month
- “Google dance”: oscillation of π during the computation

Part II

AdWords + AdSense ... or how Google turns words into money

What is it all about?

- March 2000 : the bursting of the “Internet” or “Dot-Com” Bubble
 - Many *start-ups* which offered a use value but no exchange value did not survive
 - Google had a better idea than simply selling advertising space
 - It accumulated “linguistic capital” thanks to its services
 - The idea was to exploit this capital
- An algorithm which automatically organizes speculation on words has allowed Google to create the first global linguistic market
- *Trademarks*: it was already possible to purchase certain words
- Google has boosted and liberalized that market

1

[X-HelvetiC Tours | helveticours.ch](http://helveticours.ch)www.helveticours.ch/xHelveticToursDes **vacances** à la mer aussi peu chères, ca fait vraiment du bien!

2

[Vacances tout compris | clubmed.ch](http://clubmed.ch)www.clubmed.ch/-15% sur vos **vacances** d'hiver 12/13 ou jusqu'à 480 CHF offerts now !

3

[Voyages Jusqu'à -70% - Offres Imbattables: Vos Voyages](http://groupon.ch/Voyages)www.groupon.ch/Voyages

Jusqu'à -70% avec Groupon. Ici !

[Vacances 2012. L'été à la mer, en France, Corse, Var, Bretagne ...](http://www.vacances.com/)www.vacances.com/Nombreuses annonces, de professionnels et de particuliers, pour les **vacances** d'hiver. Chalets, studios, appartements au ski. Week-ends, voyages, séjours en ...

↳ Location - Espagne - Location Aquitaine - Provence et Côte d'Azur

[Le calendrier scolaire - Ministère de l'Éducation nationale](http://www.education.gouv.fr)www.education.gouv.fr > ... > Le ministère > Repères, histoire et patrimoine**Vacances**, Zone A, Zone B, Zone C ... **Vacances** de la Toussaint ... Le départ en **vacances** a lieu après la classe, la reprise des cours le matin des jours indiqués.

↳ Vacances Scolaires 2011-2012 - Le calendrier scolaire ... - MENE0914826A

[Vacances Look Voyages : séjour pas cher en famille, club tout ...](http://www.look-voyages.fr/)www.look-voyages.fr/Votre séjour en club de **vacances** Lookea, en hôtel ou en formule circuit au meilleur prix avec Look Voyages. Départ dernière minute, pas cher ou en promo ![Villages et Clubs de vacances en tout inclus Thomas Cook](http://tt.thomascook.fr/village-club-vacances/)tt.thomascook.fr/village-club-vacances/

Ambiance. Une ambiance animée en journée comme en soirée. CHAQUE SEMAINE, NOS ANIMATEURS CONCOCTENT LE PROGRAMME DU VILLAGE : ...

[Tech Travel](http://www.techtravel.ch/)www.techtravel.ch/

les meilleures offres de voyages Culturels, Escapades, Plages

[Site Officiel Air Transat](http://www.airtransat.ch/)www.airtransat.ch/

Les plus bas prix pour les vols vers le Canada. Réservez en ligne!

[Voyager Moins Cher](http://www.voyagermoinscher.com/)www.voyagermoinscher.com/

Voyage et billet d'avion

Promotions et dernière minute.

[Vacances en France](http://www.declifrance.com/Vacance_France)www.declifrance.com/Vacance_FranceTrouvez vos **vacances** en France parmi plus de 5000 offres ![Vacances Go Voyages](http://www.govoyages.com/Vacances)www.govoyages.com/VacancesChoisissez vos **vacances** au prix le plus bas chez Go Voyages![Vacances à l'étranger](http://www.sejour-express.com/)www.sejour-express.com/Trouvez les moins chères sur le comparateur de **vacances**[Vacances en Club](http://www.club-vacances-express.fr/)www.club-vacances-express.fr/

Trouvez les moins chères sur le comparateur de location en Club

4

5

6

7

8

9

10

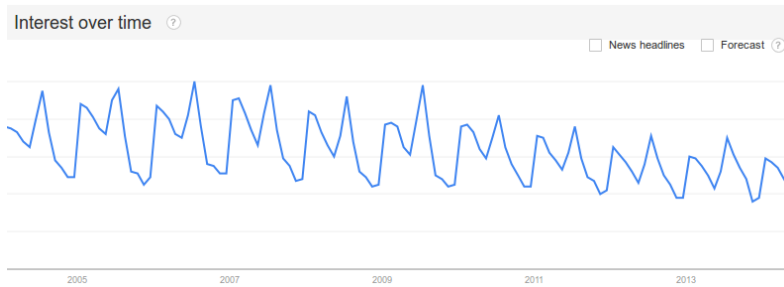
AdWords

- Auction mechanism on words to place advertisements
- All (key)words can bring about an auction
- The algorithm automatically ranks the advertisements according to a calculation in four steps:
 - Bid on a word (E): the advertiser fixes a maximum price she is willing to pay per click
 - Compute the quality score Q for the ad (relevance): **secret !**
 - Compute the rating of the ad, $R = E Q$, and its rank i
 - Compute the price to pay per click:

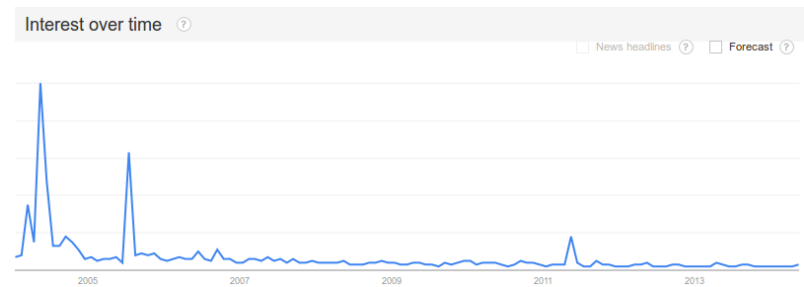
$$P_i = E_{i+1} \frac{Q_i}{Q_{i+1}}$$

Google Trends

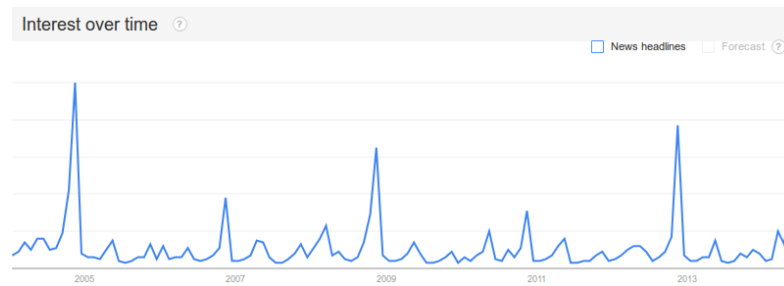
Holidays



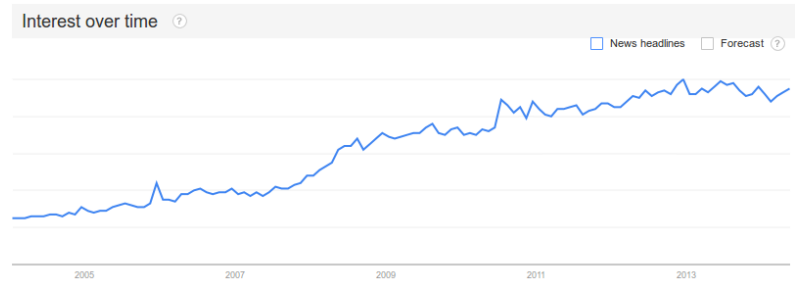
al-Qaeda



Elections



Porn



Buying and Selling Traffic



Advantages for the users

- “Free” services (search, docs, email, maps, translate, etc.)
- Useful, relevant, non-invasive advertisement
- Great user experience of on-line contents

Two Sources of Revenue

Google bicycle shops in mountain view Search About 1,820,000 results (0.33 seconds)

Local business results for **bicycle shops near Mountain View, CA**

- Performance Bicycle Shop** - www.performancebike.com
2124 West El Camino Real, Mountain View - (650) 964-1796
9 reviews, directions, and more »
- El Camino Bicycle Shop Jack'Ssee Off Ramp The** - offrampbikes.com
2320 West El Camino Real, Mountain View - (650) 968-2974
"Their customer service has always been great."
★★★★☆ 21 reviews, directions, and more »
- REI - Mountain View** - www.rei.com
2450 Charleston Road, Mountain View

Sponsored links

- Local Business Listings**
Find Phone Numbers, Links, Prices, Maps & More - Faster on Bing™!
www.Bing.com/Local
- Mountain Bikes Store**
Expert Bike Store in Your Vicinity. Locate a REI Store & Visit Us Now!
www.REI.com/BikeYourDrive
2450 Charleston Rd., Mountain View, CA
- New Bikes Up To 60% Off**
Brand Name MTBs w Full Warranties Buy Direct. Save Big. Free Shipping
www.BikesDirect.com

Advertisement on the Google sites, such as

- google.com
- gmail.com
- orkut.com
- Youtube.com



TechCrunch Gadgets Mobile Enterprise More

Google Custom Search Search

THE 88 NOW 50% SOLD!
Urban living begins in the \$300,000s
GO TO THE885J.COM San Jose, CA

About Advertise Archives Company Index Contact Events Jobs Trends

Subscribe: [Social media icons]

Advertisement on the adSense customer sites

