



Web search engines



Information Retrieval

Text Preprocessing

Inverted Index

Answering Keyword Queries

PageRank

Conclusion





Problem

*How to **index** Web content so as to answer (keyword-based) queries **efficiently**?*

Context: set of **text documents**

- d_1 The jaguar is a New World mammal of the Felidae family.
- d_2 Jaguar has designed four new engines.
- d_3 For Jaguar, Atari was keen to use a 68K family device.
- d_4 The Jacksonville Jaguars are a professional US football team.
- d_5 Mac OS X Jaguar is available at a price of US \$199 for Apple's new "family pack".
- d_6 One such ruling family to incorporate the jaguar into their name is Jaguar Paw.
- d_7 It is a big cat.

Information Retrieval

Text Preprocessing

Inverted Index

Answering Keyword Queries

PageRank

Conclusion





Text Preprocessing

Initial text preprocessing steps

- Number of optional steps
- Highly depends on the application
- Highly depends on the document language (illustrated with English)

Principle

Separate text into **tokens** (words)

Not so easy!

- In some languages (Chinese, Japanese), words **not separated by whitespace**
- Deal **consistently** with acronyms, elisions, numbers, units, URLs, emails, etc.
- **Compound words:** *hostname*, *host-name* and *host name*. Break into two tokens or regroup them as one token? In any case, lexicon and linguistic analysis needed! Even more so in other languages as German.

Usually, remove punctuation and normalize case at this point

15 février 2013



Tokenization: Example

- d_1 the₁ jaguar₂ is₃ a₄ new₅ world₆ mammal₇ of₈ the₉ felidae₁₀ family₁₁
- d_2 jaguar₁ has₂ designed₃ four₄ new₅ engines₆
- d_3 for₁ jaguar₂ atari₃ was₄ keen₅ to₆ use₇ a₈ 68k₉ family₁₀ device₁₁
- d_4 the₁ jacksonville₂ jaguars₃ are₄ a₅ professional₆ us₇ football₈ team₉
- d_5 mac₁ os₂ x₃ jaguar₄ is₅ available₆ at₇ a₈ price₉ of₁₀ us₁₁ \$199₁₂
for₁₃ apple's₁₄ new₁₅ family₁₆ pack₁₇
- d_6 one₁ such₂ ruling₃ family₄ to₅ incorporate₆ the₇ jaguar₈ into₉
their₁₀ name₁₁ is₁₂ jaguar₁₃ paw₁₄
- d_7 it₁ is₂ a₃ big₄ cat₅

Principle

Merge different forms of the same word, or of closely related words, into a single stem

- Not in all applications!
- Useful for retrieving documents containing *geese* when searching for *goose*
- Various degrees of stemming
- Possibility of building different indexes, with different stemming





Stemming schemes (1/2)

Morphological stemming (lemmatization).

- Remove **bound morphemes** from words:
 - plural markers
 - gender markers
 - tense or mood inflections
 - etc.
- Can be linguistically **very complex**, cf:
Les poules du couvent couvent.
[The hens of the monastery brood.]
- In English, somewhat **easy**:
 - Remove final -s, -'s, -ed, -ing, -er, -est
 - Take care of semiregular forms (e.g., -y/-ies)
 - Take care of irregular forms (mouse/mice)
- But still some **ambiguities**: cf rose

15 février 2013

TELECOM
ParisTech





Stemming schemes (2/2)

Lexical stemming.

- Merge **lexically related** terms of various parts of speech, such as *policy*, *politics*, *political* or *politician*
- For English, **Porter's stemming** [Porter, 1980]; stem *university* and *universal* to *univers*: not perfect!
- Possibility of coupling this with **lexicons** to merge (near-)synonyms

Phonetic stemming.

- Merge **phonetically related** words: search proper names with different spellings!
- For English, **Soundex** [US National Archives and Records Administration, 2007] stems *Robert* and *Rupert* to *R163*. Very **coarse**!

15 février 2013

TELECOM
ParisTech





Stemming Example

- d_1 the₁ jaguar₂ **be**₃ a₄ new₅ world₆ mammal₇ of₈ the₉ felidae₁₀ family₁₁
- d_2 jaguar₁ **have**₂ **design**₃ four₄ new₅ **engine**₆
- d_3 for₁ jaguar₂ atari₃ **be**₄ keen₅ to₆ use₇ a₈ 68k₉ family₁₀ device₁₁
- d_4 the₁ jacksonville₂ **jaguar**₃ **be**₄ a₅ professional₆ us₇ football₈ team₉
- d_5 mac₁ os₂ x₃ jaguar₄ **be**₅ available₆ at₇ a₈ price₉ of₁₀ us₁₁ \$199₁₂
for₁₃ **apple**₁₄ new₁₅ family₁₆ pack₁₇
- d_6 one₁ such₂ **rule**₃ family₄ to₅ incorporate₆ the₇ jaguar₈ into₉
their₁₀ name₁₁ **be**₁₂ jaguar₁₃ paw₁₄
- d_7 it₁ **be**₂ a₃ big₄ cat₅



Stop Word Removal

Principle

Remove **uninformative** words from documents, in particular to lower the cost of storing the index

determiners: *a, the, this, etc.*

function verbs: *be, have, make, etc.*

conjunctions: *that, and, etc.*

etc.





Stop Word Removal Example

- d_1 jaguar₂ new₅ world₆ mammal₇ felidae₁₀ family₁₁
- d_2 jaguar₁ design₃ four₄ new₅ engine₆
- d_3 jaguar₂ atari₃ keen₅ 68k₉ family₁₀ device₁₁
- d_4 jacksonville₂ jaguar₃ professional₆ us₇ football₈ team₉
- d_5 mac₁ os₂ x₃ jaguar₄ available₆ price₉ us₁₁ \$199₁₂ apple₁₄
new₁₅ family₁₆ pack₁₇
- d_6 one₁ such₂ rule₃ family₄ incorporate₆ jaguar₈ their₁₀ name₁₁
jaguar₁₃ paw₁₄
- d_7 big₄ cat₅



Information Retrieval

Text Preprocessing

Inverted Index

Answering Keyword Queries

PageRank

Conclusion





Inverted Index

After all preprocessing, construction of an **inverted index**:

- Index of **all terms**, with the list of documents where this term occurs
- Small scale: disk storage, with **memory mapping** (cf. `mmap`) techniques; secondary index for offset of each term in main index
- Large scale: distributed on a **cluster of machines**; hashing gives the machine responsible for a given term
- Updating the index costly, so only **batch operations** (not one-by-one addition of term occurrences)





Inverted Index Example

family	d_1, d_3, d_5, d_6
football	d_4
jaguar	$d_1, d_2, d_3, d_4, d_5, d_6$
new	d_1, d_2, d_5
rule	d_6
us	d_4, d_5
world	d_1
...	



Storing positions in the index

- phrase queries, NEAR operator: need to keep position information in the index
- just add it in the document list!

family $d_1/11, d_3/10, d_5/16, d_6/4$

football $d_4/8$

jaguar $d_1/2, d_2/1, d_3/2, d_4/3, d_5/4, d_6/8 + 13$

new $d_1/5, d_2/5, d_5/15$

rule $d_6/3$

us $d_4/7, d_5/11$

world $d_1/6$

...





TF-IDF Weighting

- Some term occurrences have more **weight** than others:
 - Terms occurring **frequently** in a **given document**: more **relevant**
 - Terms occurring **rarely** in the **document collection** as a whole: more **informative**
- Add **Term Frequency—Inverse Document Frequency** weighting to occurrences;

$$\text{tfidf}(t, d) = \frac{n_{t,d}}{\sum_{t'} n_{t',d}} \cdot \log \frac{|D|}{|\{d' \in D \mid n_{t,d'} > 0\}|}$$

$n_{t,d}$ number of occurrences of t in d
 D set of all documents

- Store documents (along with weight) in **decreasing weight order** in the index



TF-IDF Weighting Example

family	$d_1/11/.13, d_3/10/.13, d_6/4/.08, d_5/16/.07$
football	$d_4/8/.47$
jaguar	$d_1/2/.04, d_2/1/.04, d_3/2/.04, d_4/3/.04, d_6/8 + 13/.04, d_5/4/.$
new	$d_2/5/.24, d_1/5/.20, d_5/15/.10$
rule	$d_6/3/.28$
us	$d_4/7/.30, d_5/11/.15$
world	$d_1/6/.47$
...	



Information Retrieval

Text Preprocessing

Inverted Index

Answering Keyword Queries

PageRank

Conclusion





Answering Boolean Queries

- **Single keyword query**: just consult the index and return the documents in index order.
- **Boolean multi-keyword query**

(jaguar AND new AND NOT family) OR cat

Same way! Retrieve document lists from all keywords and apply adequate set operations:

AND intersection

OR union

AND NOT difference

- **Global score**: some function of the individual weight (e.g., addition for conjunctive queries)
- **Position queries**: consult the index, and filter by appropriate condition

15 février 2013



Answering Top- k Queries

$t_1 \text{ AND } \dots \text{ AND } t_n$

$t_1 \text{ OR } \dots \text{ OR } t_n$

Problem

Find the **top- k results** (for some given k) to the query, without retrieving all documents matching it.

Notations:

$s(t, d)$ weight of t in d (e.g., tfidf)

$g(s_1, \dots, s_n)$ monotonous function that computes the global score
(e.g., addition)





Fagin's Threshold Algorithm

[Fagin et al., 2001]

(with an additional direct index giving $s(t, d)$)

1. Let R be the empty list and $m = +\infty$.
2. For each $1 \leq i \leq n$:
 - 2.1 Retrieve the document $d^{(i)}$ containing term t_i that has the **next largest** $s(t_i, d^{(i)})$.
 - 2.2 Compute its global score $g_{d^{(i)}} = g(s(t_1, d^{(i)}), \dots, s(t_n, d^{(i)}))$ by retrieving all $s(t_j, d^{(i)})$ with $j \neq i$.
 - 2.3 If R contains less than k documents, or if $g_{d^{(i)}}$ is greater than the **minimum of the score of documents** in R , add $d^{(i)}$ to R (and remove the worst element in R if it is full).
3. Let $m = g(s(t_1, d^{(1)}), s(t_2, d^{(2)}), \dots, s(t_n, d^{(n)}))$.
4. If R contains k documents, and the minimum of the score of the documents in R is **greater than or equal** to m , return R .
5. Redo step 2.



Information Retrieval

PageRank

Conclusion





The Web Graph

The World Wide Web seen as a (directed) graph:

Vertices: Web pages

Edges: hyperlinks

Same for other interlinked environments:

- dictionaries
- encyclopedias
- scientific publications
- social networks





PageRank (Google's Ranking [Brin and Page, 1998])

Idea

Important pages are pages pointed to by **important** pages.

$$\begin{cases} g_{ij} = 0 & \text{if there is no link between page } i \text{ and } j; \\ g_{ij} = \frac{1}{n_i} & \text{otherwise, with } n_i \text{ the number of outgoing links of page } i. \end{cases}$$

Definition (Tentative)

Probability that the surfer following the **random walk** in G has arrived on page i at some distant given point in the future.

$$\text{pr}(i) = \left(\lim_{k \rightarrow +\infty} (G^T)^k v \right)_i$$

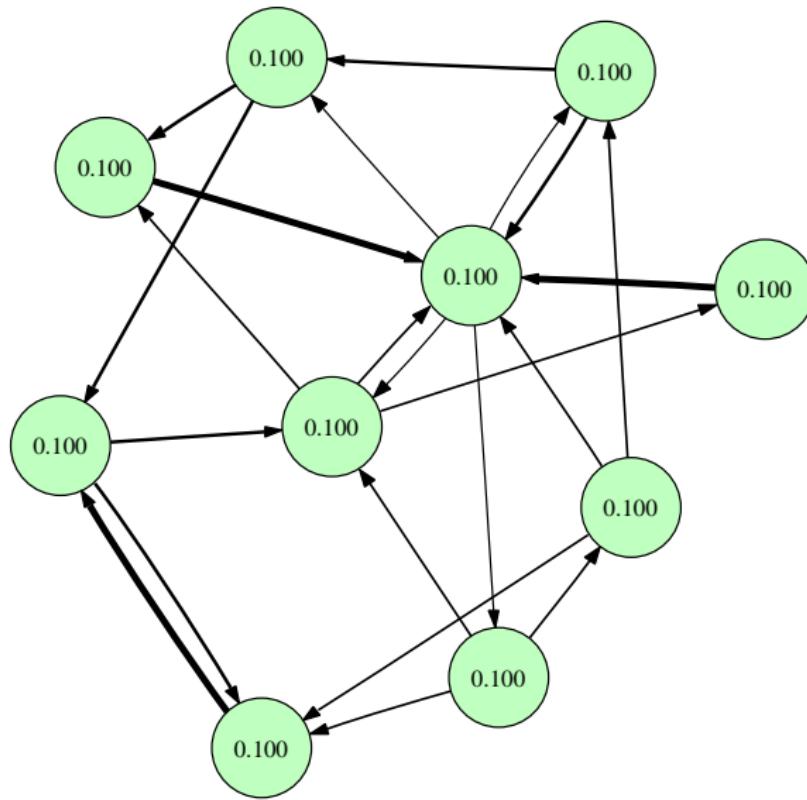
where v is some initial column vector.

15 février 2013



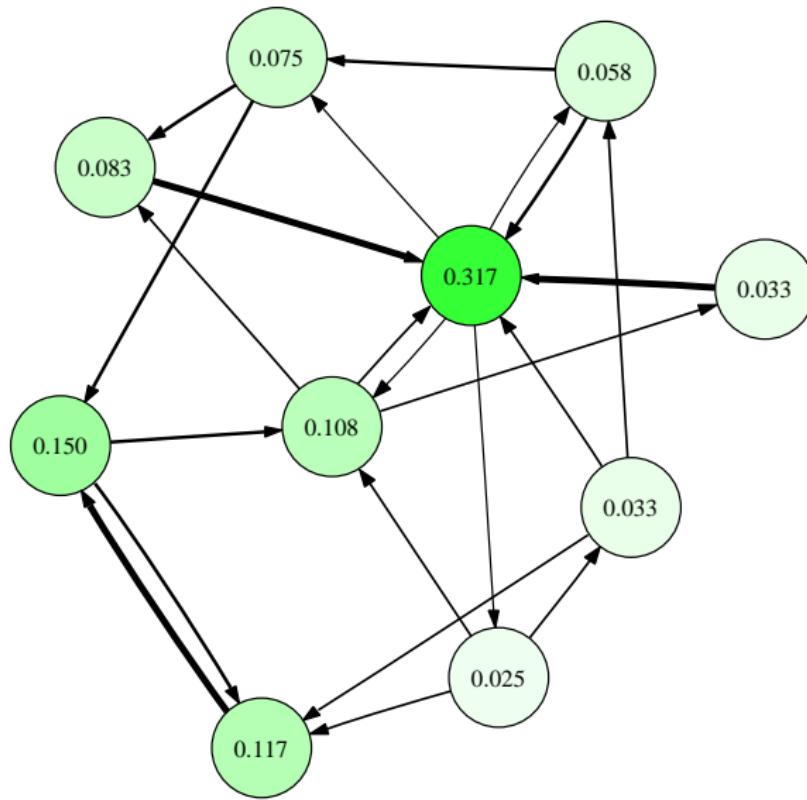


Illustrating PageRank Computation



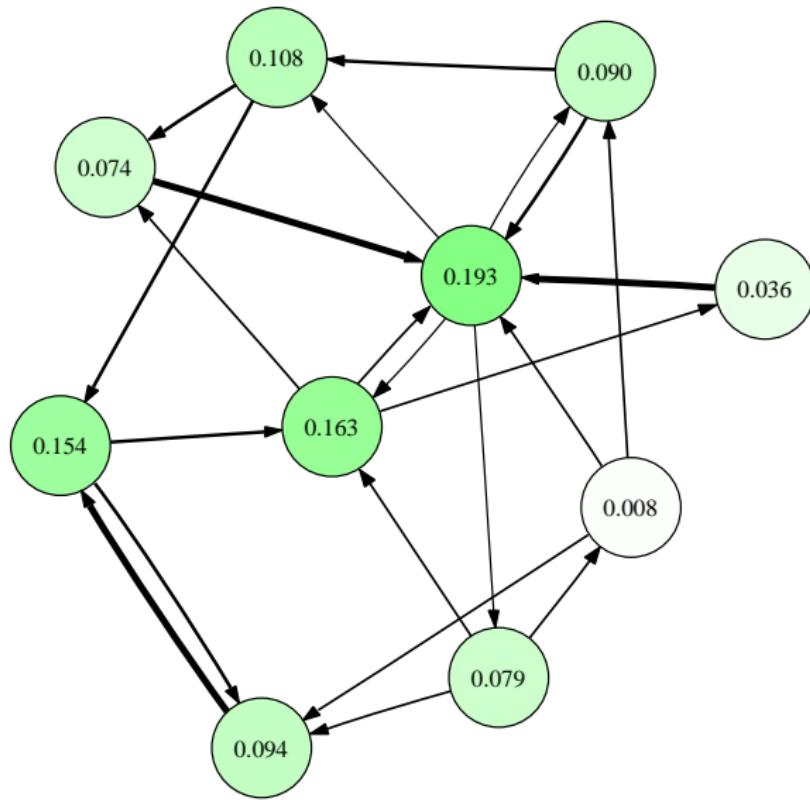


Illustrating PageRank Computation



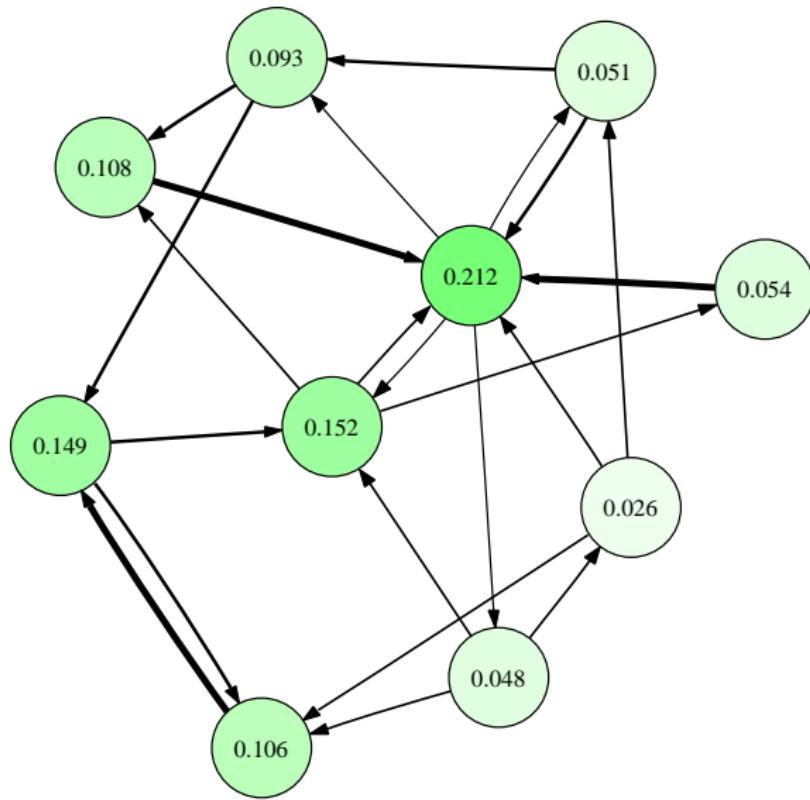


Illustrating PageRank Computation



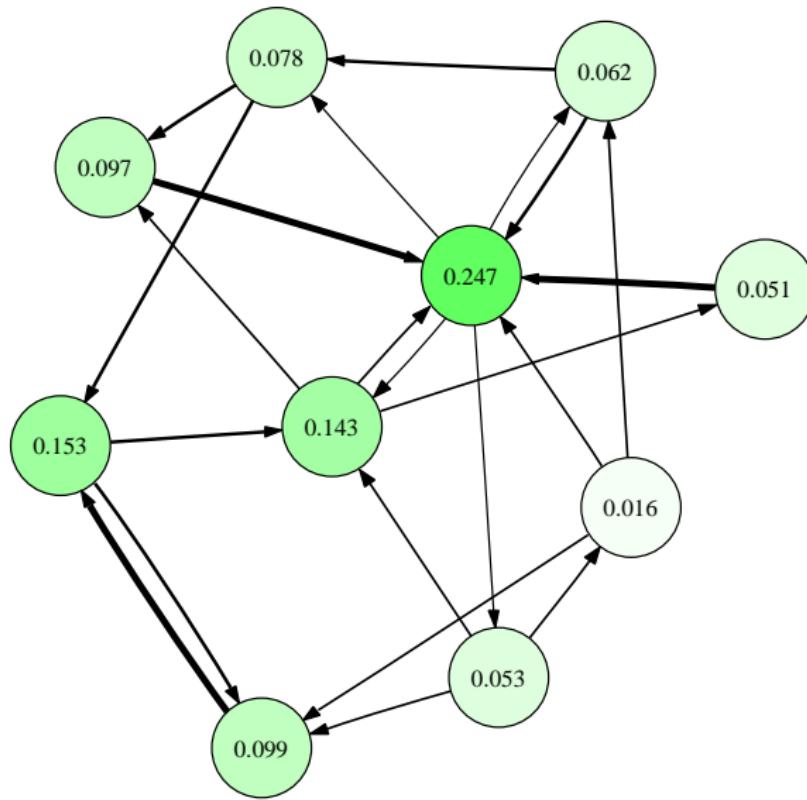


Illustrating PageRank Computation



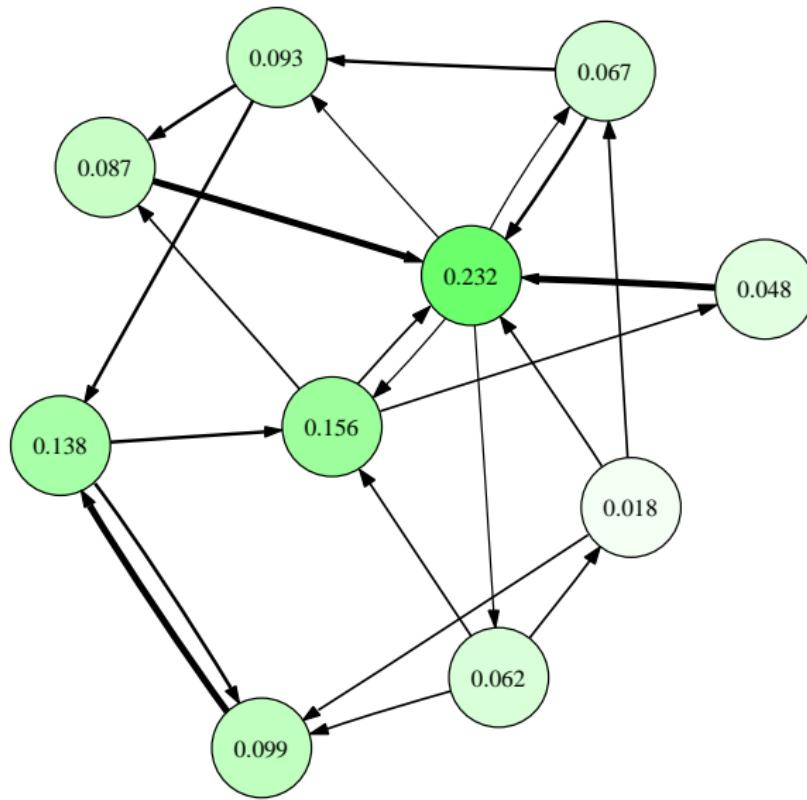


Illustrating PageRank Computation



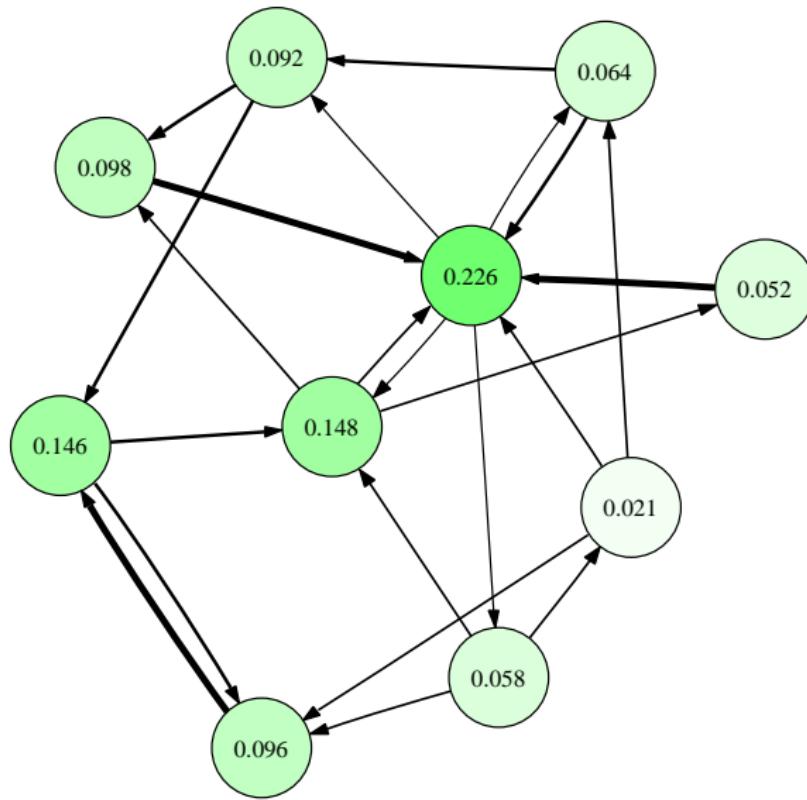


Illustrating PageRank Computation



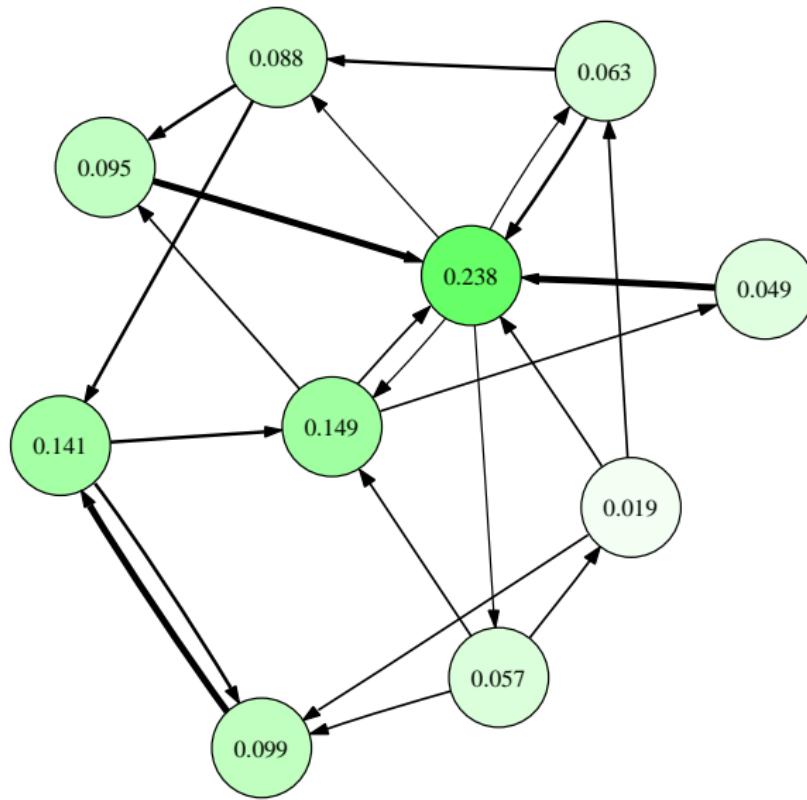


Illustrating PageRank Computation



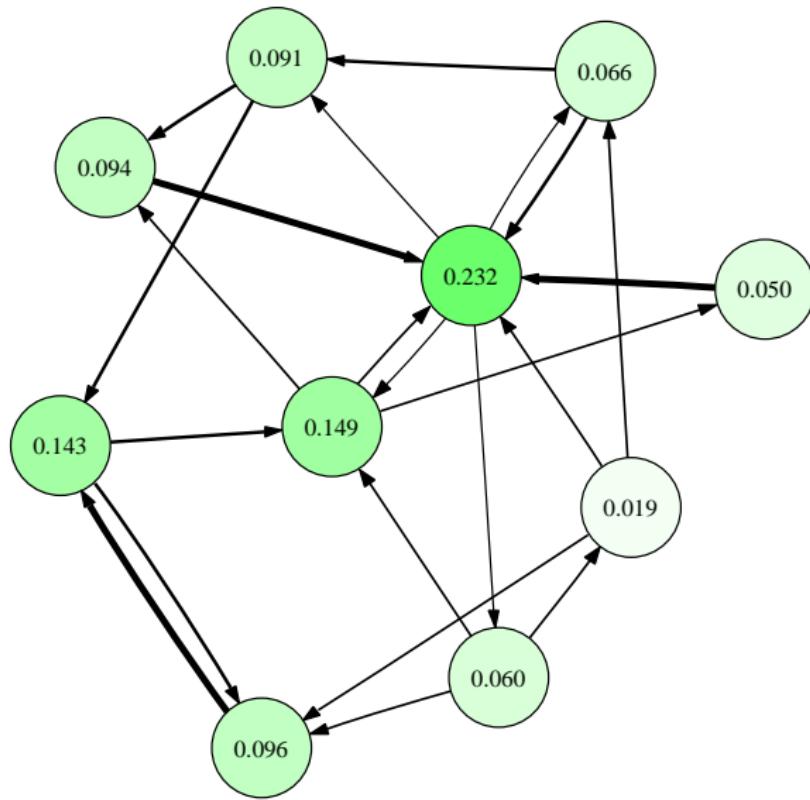


Illustrating PageRank Computation



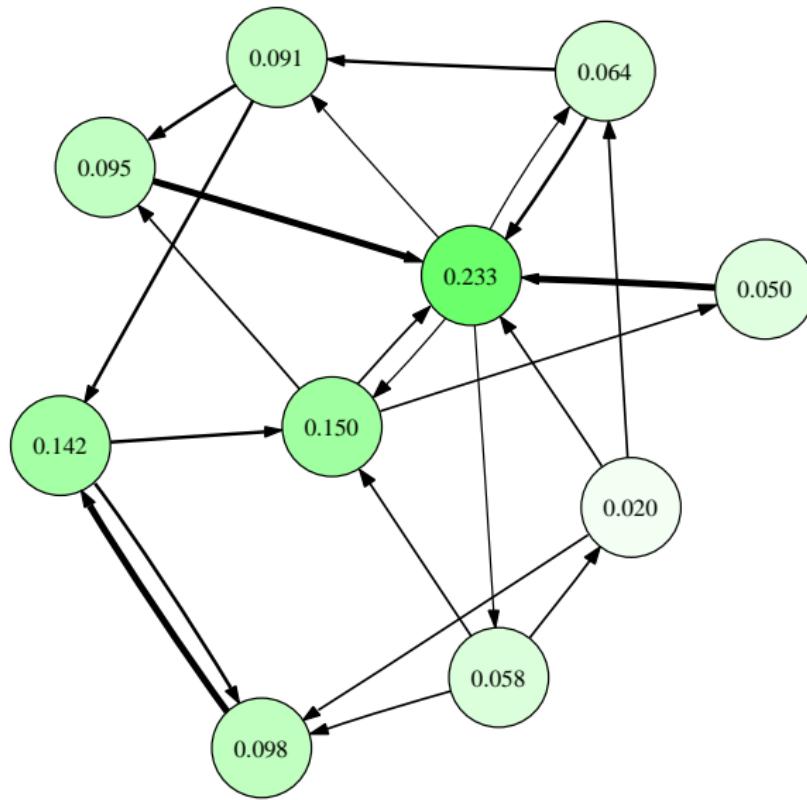


Illustrating PageRank Computation



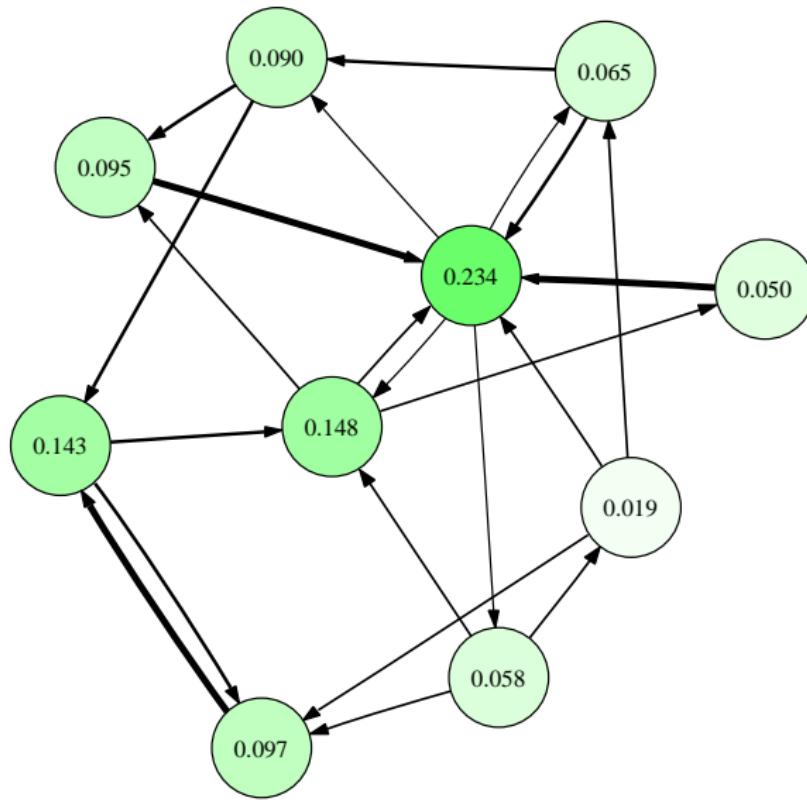


Illustrating PageRank Computation



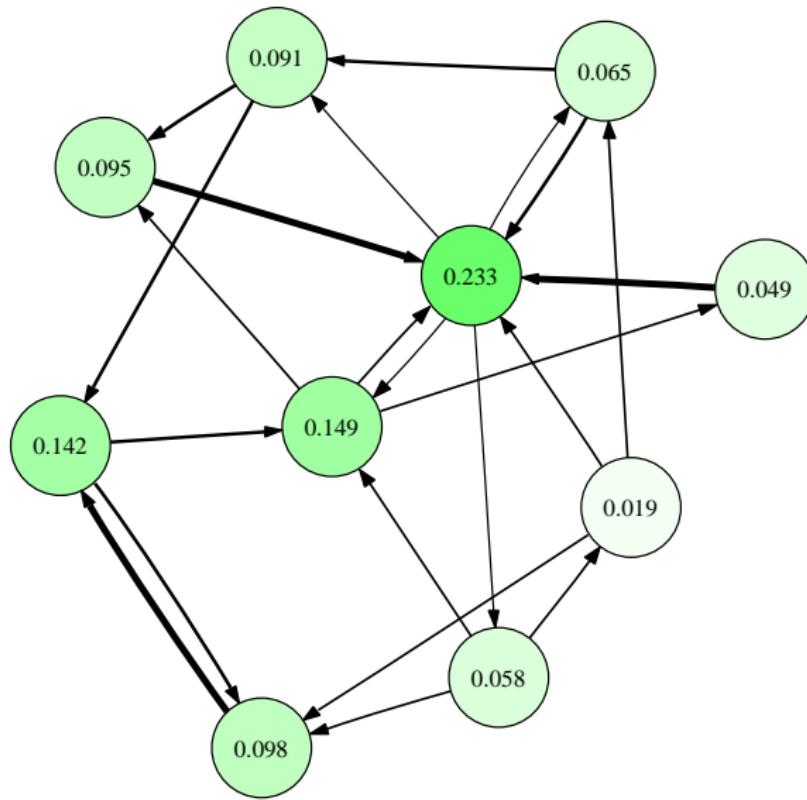


Illustrating PageRank Computation



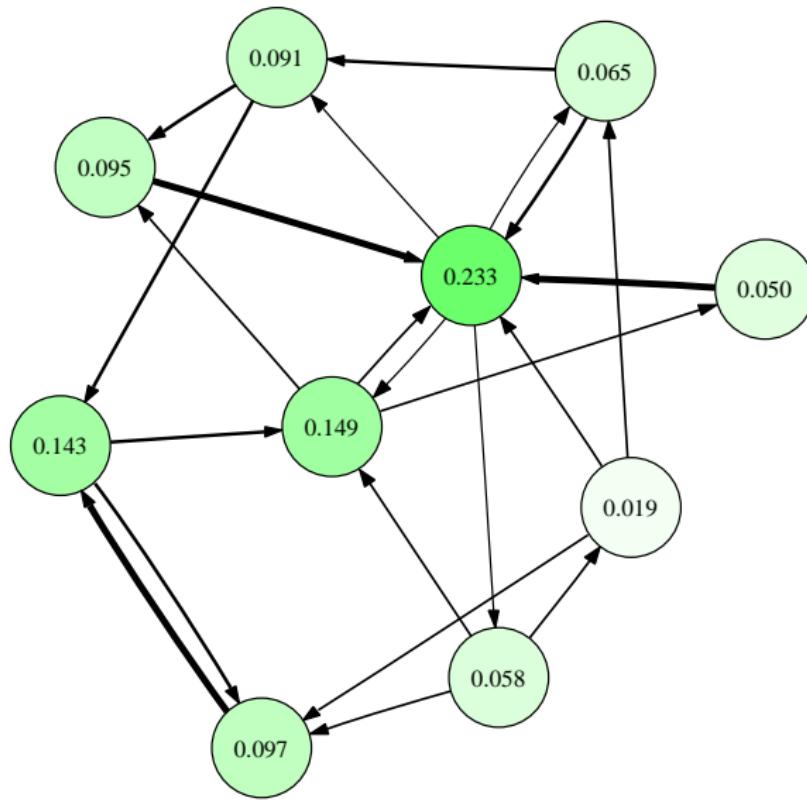


Illustrating PageRank Computation



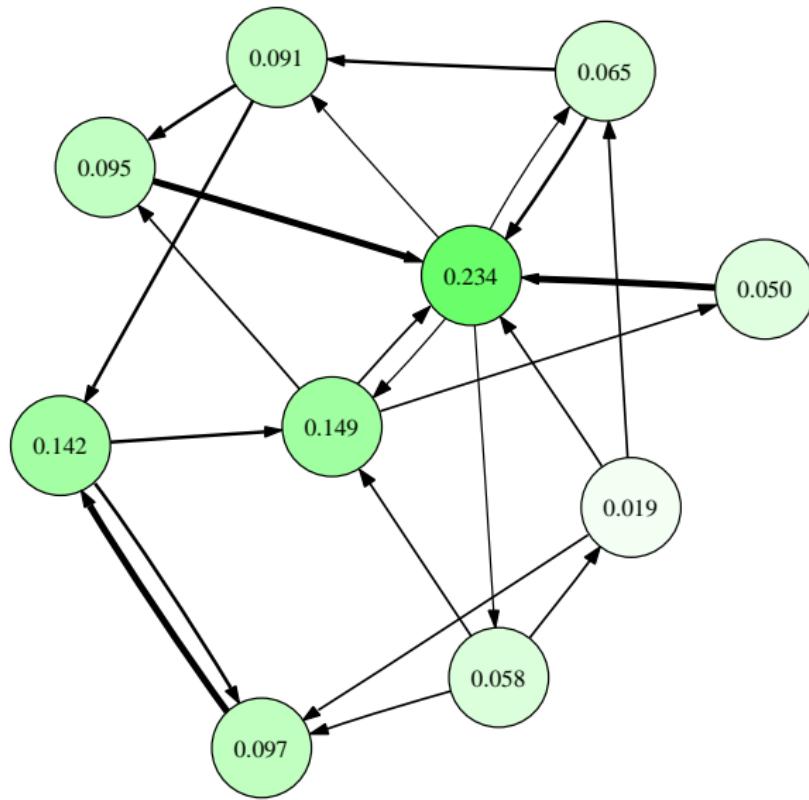


Illustrating PageRank Computation





Illustrating PageRank Computation





PageRank with Damping

May not always converge, or convergence may not be unique.

To fix this, the random surfer can at each step randomly jump to any page of the Web with some probability d ($1 - d$: damping factor).

$$\text{pr}(i) = \left(\lim_{k \rightarrow +\infty} ((1 - d)G^T + dU)^k v \right)_i$$

where U is the matrix with all $\frac{1}{N}$ values with N the number of vertices.





Using PageRank to Score Query Results

- PageRank: **global** score, independent of the query
- Can be used to raise the weight of **important** pages:

$$\text{weight}(t, d) = \text{tfidf}(t, d) \times \text{pr}(d),$$

- This can be directly incorporated **in the index**.





Iterative Computation of PageRank

1. Compute G (often stored as its adjacency list). Make sure lines sum to 1.
2. Let u be the uniform vector of sum 1, $v = u$, w the zero vector.
3. While v is **different enough** from w :
 - Set $w = v$.
 - Set $v = (1 - d)G^T v + du$.

Information Retrieval

PageRank

Conclusion



What you should remember

- The inverted index model, associated tools and techniques
- Main ideas behind Fagin's TA
- PageRank, and its iterative computation

Software

- Most DBMS have text indexing capabilities (e.g., MySQL's FULLTEXT indexes)
- Apache Lucene, free software library for information retrieval

To go further

- Good textbooks [Manning et al., 2008, Chakrabarti, 2003]
- Very influential papers on top- k algorithms and PageRank: [Fagin et al., 2001, Brin and Page, 1998]

Bibliography I

- Sergey Brin and Lawrence Page. The anatomy of a large-scale hypertextual Web search engine. *Computer Networks*, 30(1–7): 107–117, April 1998.
- Soumen Chakrabarti. *Mining the Web: Discovering Knowledge from Hypertext Data*. Morgan Kaufmann, San Francisco, USA, 2003.
- Ronald Fagin, Amnon Lotem, and Moni Naor. Optimal aggregation algorithms for middleware. In *Proc. PODS*, Santa Barbara, USA, May 2001.
- Christopher D. Manning, Prabhakar Raghavan, and Hinrich Schütze. *Introduction to Information Retrieval*. Cambridge University Press, Cambridge, United Kingdom, 2008. Available online at <http://informationretrieval.org>.
- Martin F. Porter. An algorithm for suffix stripping. *Program*, 14(3): 130–137, July 1980.

Bibliography II

US National Archives and Records Administration. The Soundex indexing system.

<http://www.archives.gov/genealogy/census/soundex.html>, May 2007.



Licence de droits d'usage



Contexte public } avec modifications

Par le téléchargement ou la consultation de ce document, l'utilisateur accepte la licence d'utilisation qui y est attachée, telle que détaillée dans les dispositions suivantes, et s'engage à la respecter intégralement.

La licence confère à l'utilisateur un droit d'usage sur le document consulté ou téléchargé, totalement ou en partie, dans les conditions définies ci-après et à l'exclusion expresse de toute utilisation commerciale.

Le droit d'usage défini par la licence autorise un usage à destination de tout public qui comprend :

- le droit de reproduire tout ou partie du document sur support informatique ou papier,
- le droit de diffuser tout ou partie du document au public sur support papier ou informatique, y compris par la mise à la disposition du public sur un réseau numérique,
- le droit de modifier la forme ou la présentation du document,
- le droit d'intégrer tout ou partie du document dans un document composite et de le diffuser dans ce nouveau document, à condition que :
 - L'auteur soit informé.

Les mentions relatives à la source du document et/ou à son auteur doivent être conservées dans leur intégralité.

Le droit d'usage défini par la licence est personnel et non exclusif.

Tout autre usage que ceux prévus par la licence est soumis à autorisation préalable et expresse de l'auteur : sitepedago@telecom-paristech.fr

