

# Web mining

## Information Retrieval

9 October 2014





# Inverted Index Model

Text Preprocessing

Inverted Index

Answering Keyword Queries

Indexing Other Media

PageRank

Search Engine Optimization

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.





Inverted Index Model

Text Preprocessing

Inverted Index

Answering Keyword Queries

Indexing Other Media

PageRank

Search Engine Optimization

Conclusion





## Initial text preprocessing steps

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





How to find the language used in a document?

- Meta-information about the document: often **not reliable!**
- **Unambiguous** scripts or letters: not very common!

한글

カタカナ

مِرقِ

Gharbi

porn





How to find the language used in a document?

- Meta-information about the document: often **not reliable!**
- **Unambiguous** scripts or letters: not very common!

한글

カタカナ

ދިވެހި

Għarbi

þorn

**Respectively:** Korean Hanguk, Japanese Katakana, Maldivian Dhivehi, Maltese, Icelandic

- Extension of this: **frequent characters**, or, better, **frequent k-grams**
- Use standard machine learning techniques (**classifiers**)



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





- d*<sub>1</sub> the<sub>1</sub> jaguar<sub>2</sub> is<sub>3</sub> a<sub>4</sub> new<sub>5</sub> world<sub>6</sub> mammal<sub>7</sub> of<sub>8</sub> the<sub>9</sub> felidae<sub>10</sub> family<sub>11</sub>
- d*<sub>2</sub> jaguar<sub>1</sub> has<sub>2</sub> designed<sub>3</sub> four<sub>4</sub> new<sub>5</sub> engines<sub>6</sub>
- d*<sub>3</sub> for<sub>1</sub> jaguar<sub>2</sub> atari<sub>3</sub> was<sub>4</sub> keen<sub>5</sub> to<sub>6</sub> use<sub>7</sub> a<sub>8</sub> 68k<sub>9</sub> family<sub>10</sub> device<sub>11</sub>
- d*<sub>4</sub> the<sub>1</sub> jacksonville<sub>2</sub> jaguars<sub>3</sub> are<sub>4</sub> a<sub>5</sub> professional<sub>6</sub> us<sub>7</sub> football<sub>8</sub> team<sub>9</sub>
- d*<sub>5</sub> mac<sub>1</sub> os<sub>2</sub> x<sub>3</sub> jaguar<sub>4</sub> is<sub>5</sub> available<sub>6</sub> at<sub>7</sub> a<sub>8</sub> price<sub>9</sub> of<sub>10</sub> us<sub>11</sub> \$199<sub>12</sub>  
for<sub>13</sub> apple's<sub>14</sub> new<sub>15</sub> family<sub>16</sub> pack<sub>17</sub>
- d*<sub>6</sub> one<sub>1</sub> such<sub>2</sub> ruling<sub>3</sub> family<sub>4</sub> to<sub>5</sub> incorporate<sub>6</sub> the<sub>7</sub> jaguar<sub>8</sub> into<sub>9</sub>  
their<sub>10</sub> name<sub>11</sub> is<sub>12</sub> jaguar<sub>13</sub> paw<sub>14</sub>
- d*<sub>7</sub> it<sub>1</sub> is<sub>2</sub> a<sub>3</sub> big<sub>4</sub> cat<sub>5</sub>



**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

## 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



## 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**!





- d*<sub>1</sub> the<sub>1</sub> jaguar<sub>2</sub> **be**<sub>3</sub> a<sub>4</sub> new<sub>5</sub> world<sub>6</sub> mammal<sub>7</sub> of<sub>8</sub> the<sub>9</sub> felidae<sub>10</sub> family<sub>11</sub>
- d*<sub>2</sub> jaguar<sub>1</sub> **have**<sub>2</sub> **design**<sub>3</sub> four<sub>4</sub> new<sub>5</sub> **engine**<sub>6</sub>
- d*<sub>3</sub> for<sub>1</sub> jaguar<sub>2</sub> atari<sub>3</sub> **be**<sub>4</sub> keen<sub>5</sub> to<sub>6</sub> use<sub>7</sub> a<sub>8</sub> 68k<sub>9</sub> family<sub>10</sub> device<sub>11</sub>
- d*<sub>4</sub> the<sub>1</sub> jacksonville<sub>2</sub> **jaguar**<sub>3</sub> **be**<sub>4</sub> a<sub>5</sub> professional<sub>6</sub> us<sub>7</sub> football<sub>8</sub> team<sub>9</sub>
- d*<sub>5</sub> mac<sub>1</sub> os<sub>2</sub> x<sub>3</sub> jaguar<sub>4</sub> **be**<sub>5</sub> available<sub>6</sub> at<sub>7</sub> a<sub>8</sub> price<sub>9</sub> of<sub>10</sub> us<sub>11</sub> \$199<sub>12</sub>  
for<sub>13</sub> **apple**<sub>14</sub> new<sub>15</sub> family<sub>16</sub> pack<sub>17</sub>
- d*<sub>6</sub> one<sub>1</sub> such<sub>2</sub> **rule**<sub>3</sub> family<sub>4</sub> to<sub>5</sub> incorporate<sub>6</sub> the<sub>7</sub> jaguar<sub>8</sub> into<sub>9</sub>  
their<sub>10</sub> name<sub>11</sub> **be**<sub>12</sub> jaguar<sub>13</sub> paw<sub>14</sub>
- d*<sub>7</sub> it<sub>1</sub> **be**<sub>2</sub> a<sub>3</sub> big<sub>4</sub> cat<sub>5</sub>





## 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.





- $d_1$  jaguar<sub>2</sub> new<sub>5</sub> world<sub>6</sub> mammal<sub>7</sub> felidae<sub>10</sub> family<sub>11</sub>
- $d_2$  jaguar<sub>1</sub> design<sub>3</sub> four<sub>4</sub> new<sub>5</sub> engine<sub>6</sub>
- $d_3$  jaguar<sub>2</sub> atari<sub>3</sub> keen<sub>5</sub> 68k<sub>9</sub> family<sub>10</sub> device<sub>11</sub>
- $d_4$  jacksonville<sub>2</sub> jaguar<sub>3</sub> professional<sub>6</sub> us<sub>7</sub> football<sub>8</sub> team<sub>9</sub>
- $d_5$  mac<sub>1</sub> os<sub>2</sub> x<sub>3</sub> jaguar<sub>4</sub> available<sub>6</sub> price<sub>9</sub> us<sub>11</sub> \$199<sub>12</sub> apple<sub>14</sub>  
new<sub>15</sub> family<sub>16</sub> pack<sub>17</sub>
- $d_6$  one<sub>1</sub> such<sub>2</sub> rule<sub>3</sub> family<sub>4</sub> incorporate<sub>6</sub> jaguar<sub>8</sub> their<sub>10</sub> name<sub>11</sub>  
jaguar<sub>13</sub> paw<sub>14</sub>
- $d_7$  big<sub>4</sub> cat<sub>5</sub>





Inverted Index Model

Text Preprocessing

**Inverted Index**

Answering Keyword Queries

Indexing Other Media

PageRank

Search Engine Optimization

Conclusion

9 October 2014





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)





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	$d_1$
...	





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	$d_1/6$
...	





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





family  $d_1/11/.13, d_3/10/.13, d_6/4/.08, d_5/16/.07$   
football  $d_4/8/.47$   
jaguar  $d_1/2/.04, d_2/1/.04, d_3/2/.04, d_4/3/.04, d_6/8 + 13/.04, d_5/4/.04$   
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  $d_1/6/.47$   
...





# Inverted Index Model

Text Preprocessing

Inverted Index

Answering Keyword Queries

Indexing Other Media

PageRank

Search Engine Optimization

Conclusion

9 October 2014





- **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





$t_1$  AND ... AND  $t_n$

$t_1$  OR ... 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)



(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.





## The Inverted Index Model

### Indexing Other Media

### PageRank

### Search Engine Optimization

### Conclusion





- HTML: text + meta-information + structure
- Possibly: separate index for meta-information (title, keywords)
- Increase weight of structurally emphasized content in index
- Tree structure can also be queried with XPath or XQuery, but not very useful on the Web as a whole, because of tag soup and lack of consistency.





Basic approach: index **text from context** of the media

- surrounding text
  - text in or around the links pointing to the content
  - filenames
  - associated subtitles (hearing-impaired track on TV)
- Elaborate approach: index and search the media itself, with the help of **speech recognition** and **sound, image, and video analysis**.  
Mostly still experimental!
- TrackID, Shazam: identify a song played on the radio
  - Musipedia: look for music by whistling a tune,  
<http://www.musipedia.org/>, <http://www.midomi.com/>
  - Image search from a similar image, <http://images.google.com/>,  
Google Goggles
  - Voxlead, <http://voxleadnews.labs.exalead.com/>





The Inverted Index Model

Indexing Other Media

PageRank

Search Engine Optimization

Conclusion



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)

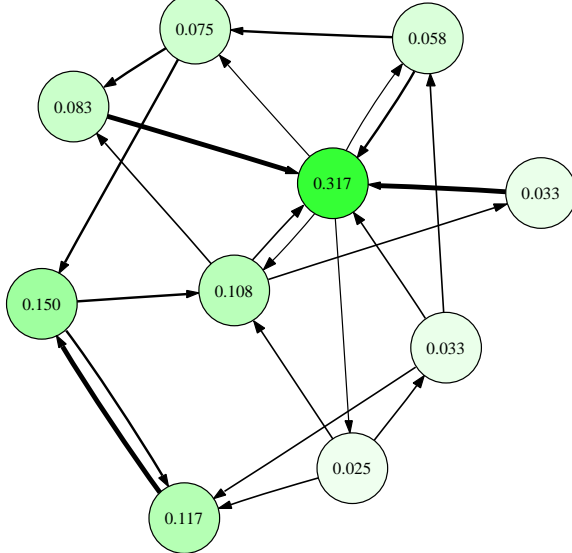
**Probability** that the surfer following the **random walk** in  $G$  has arrived on page  $i$  at some distant given point in the future.

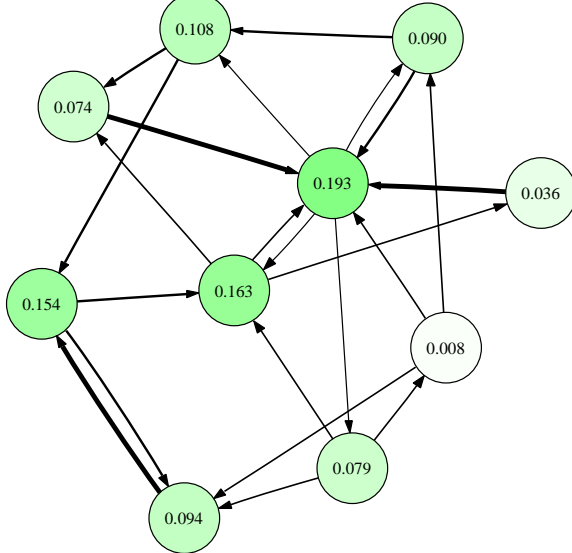
$$\text{pr}(i) = \left( \lim_{k \rightarrow +\infty} (G^T)^k v \right)_i$$

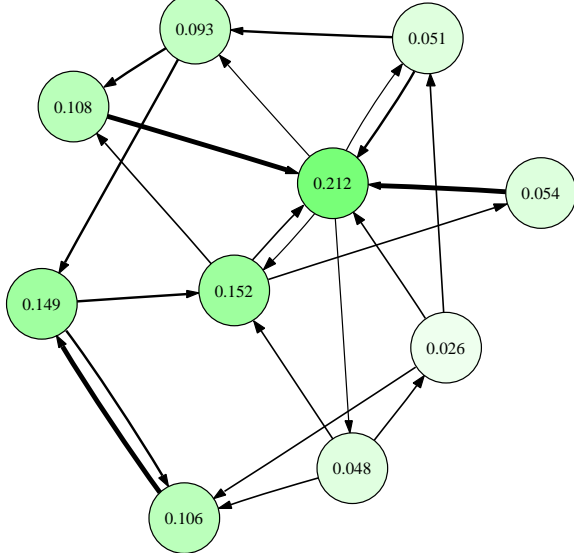
where  $v$  is some initial column vector.

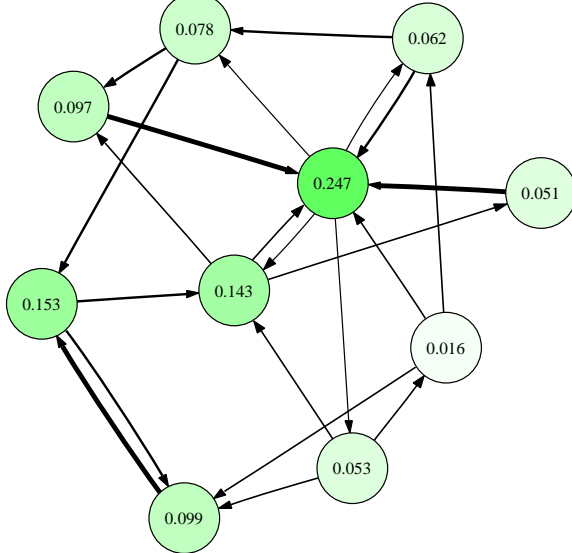


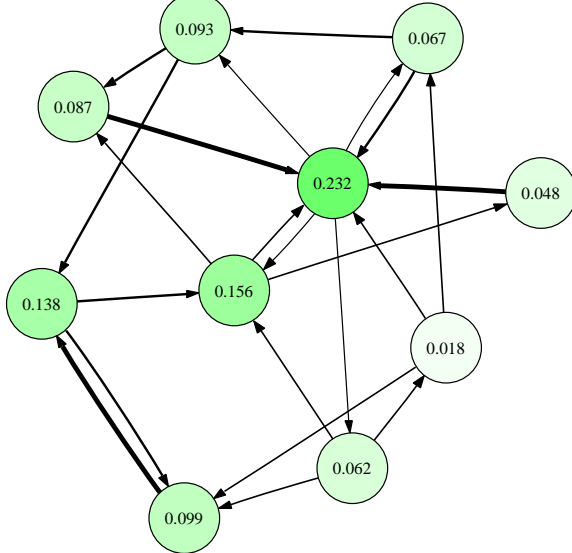


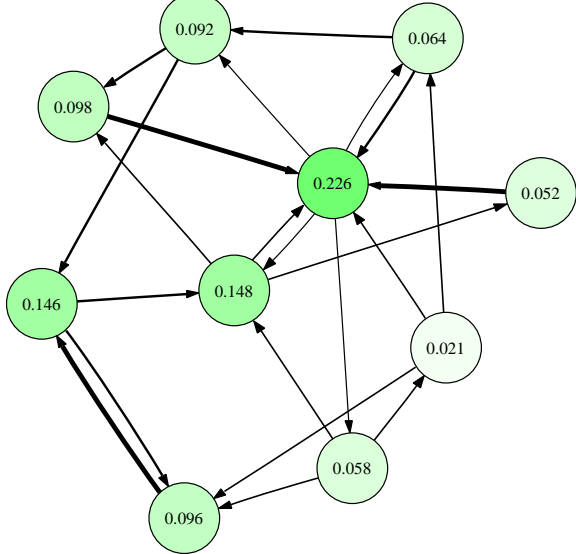


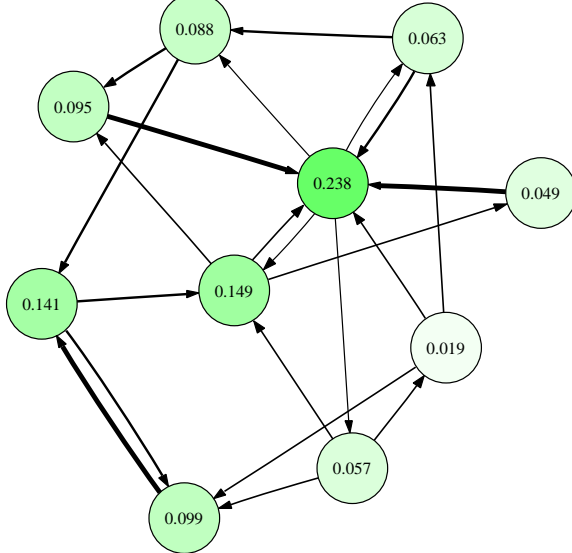


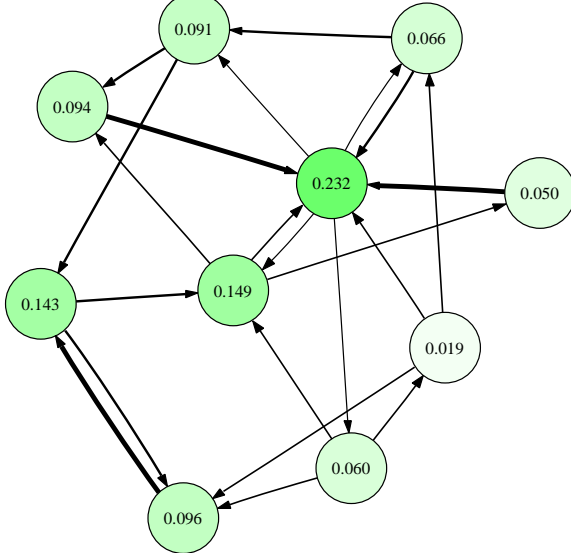


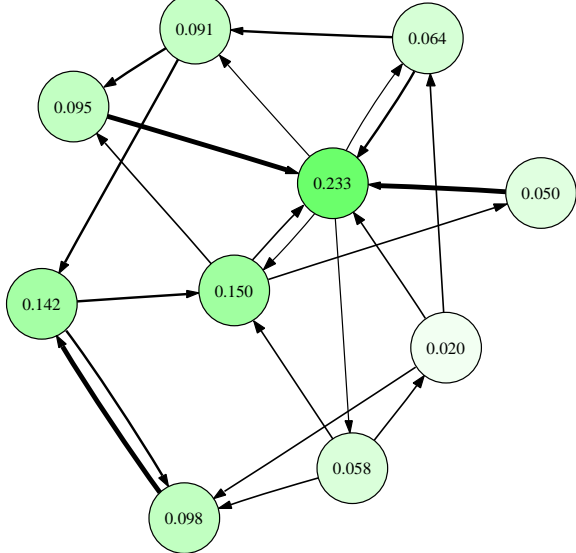


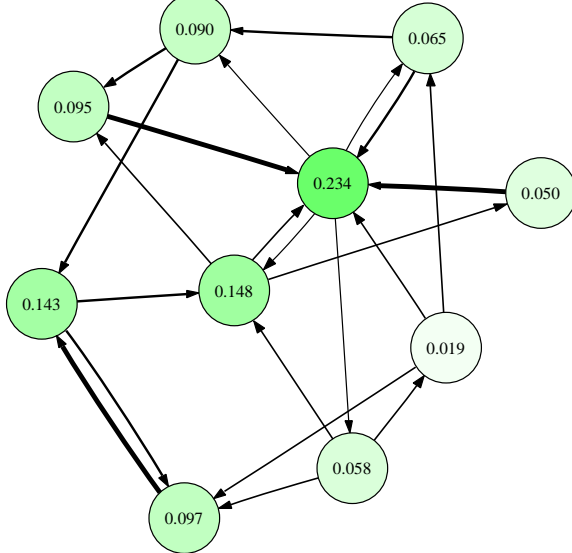


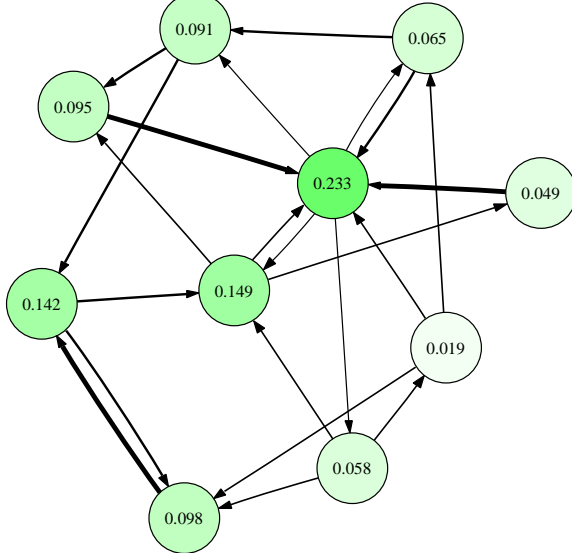


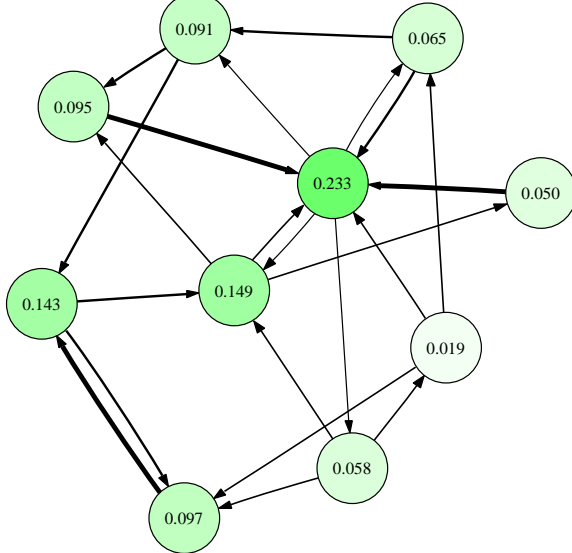


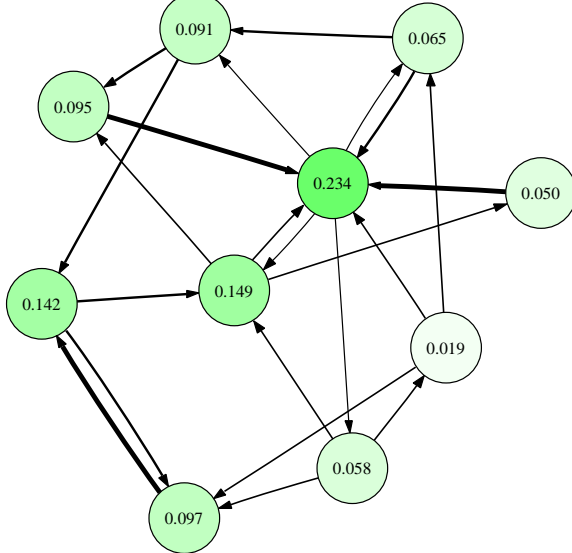














May not always converge, or convergence may not be unique.  
To fix this, the random surfer can at each step randomly jump to any page of the Web with some probability  $d$  ( $1 - d$ : damping factor).

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

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





- 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**.





The Inverted Index Model

Indexing Other Media

PageRank

Search Engine Optimization

Spamdexing

Optimization

Conclusion





The Inverted Index Model

Indexing Other Media

PageRank

Search Engine Optimization

Spamdexing

Optimization

Conclusion





## Definition

Fraudulent techniques that are used by unscrupulous webmasters to artificially raise the visibility of their website to users of search engines

Purpose: attracting visitors to websites to make profit.

Unceasing war between **spamdexers** and **search engines**



Put **unrelated** terms in:

- meta-information (<meta name="description">, <meta name="keywords">)
- text content hidden to the user with JavaScript, CSS, or HTML presentational elements

## Countertechnique

- **Ignore** meta-information
- Try and **detect** invisible text

Huge number of hosts on the Internet used for the sole purpose of **referencing** each other, without any content in themselves, to **raise the importance** of a given website or set of websites.

## Countertechnique

- Detection of websites with **empty** or **duplicate** content
- Use of heuristics to discover **subgraphs** that look like link farms



## Technique

Pollute **user-editable** websites (blogs, wikis) or exploit security bugs to add **artificial** links to websites, in order to raise its importance.

## Countertechnique

**rel="nofollow"** attribute to `<a>` links not validated by a page's owner





The Inverted Index Model

Indexing Other Media

PageRank

Search Engine Optimization

Spamdexing

Optimization

Conclusion





Faire attention à l'accessibilité aux robots des moteurs de recherche (Flash, JavaScript, etc.)

- Éventuellement prévoir des versions HTML pures alternatives
- Site accessible à une URL unique (par exemple, entre `http://www.toto.com/` et `http://toto.com/`, si les deux permettent d'accéder au site, l'un doit être une redirection (HTTP) vers l'autre)
- Structure cohérente (hébergement sur un seul serveur, contenu des URLs pertinentes), bonne structure de liens (de n'importe quelle page, il doit être possible d'atteindre la page principale en un clic, et n'importe quelle autre page en quelques clics)
- Liens vers les autres sites pertinents, que le site n'apparaisse pas complètement isolé





- Pas de tentatives de spamdexing: pas de texte invisible, pas de liens dans les mots-clefs de la balise <meta> (assez inutiles de toute façon), etc.
- Faire apparaître des liens vers le site (surtout vers la page principale, et éventuellement pages internes quand c'est pertinent) sur d'autres sites Web (référencement dans des annuaires, sur les pages Web des individus/sociétés en lien avec le contenu du site, etc.).
- Éventuellement s'il est difficile de faire ainsi, soumettre le site aux différents moteurs de recherche (p. ex., <http://www.google.com/addurl/> pour Google).
- Se méfier des sociétés proposant un meilleur référencement contre rémunération; leurs pratiques sont mal vues par les moteurs de recherche.
- Et enfin... avoir du contenu intéressant !





The Inverted Index Model

Indexing Other Media

PageRank

Search Engine Optimization

Conclusion



## What you should remember



The inverted index model, associated tools and techniques

- Main ideas behind Fagin's TA and NRA
- The document vector space model

## Software

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

## To go further

- A good textbook [Manning et al., 2008]. Available online, along with slides!
- A very influential paper on top- $k$  algorithms: [Fagin et al., 2001]



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.

US National Archives and Records Administration. The Soundex indexing system. <http://www.archives.gov/genealogy/census/soundex.html>, May 2007.



**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](mailto:sitepedago@telecom-paristech.fr)

