# Collaborative(ly) Personalized PageRank

September 2014
Théo Dubourg - theo.dubourg@insa-lyon.fr

# Hey ! I am a slide.

I am projected, shiny, I'm quite big, but it is not mandatory to read me.

Listening is enough, if needed to look at me, the guy over there will tell you.

# Context: Web Search, IR Limits

- Only use the text/document

- Not aware of the "environment"

- Human beings are not really compatible with IR: They do not express themselves the right way.
  "Human beings do not know how to search" - Someone

# Context: Web Search, Some Solutions

Use the environment/specificities of the web:

- Linking information: PageRank

- Techniques that Google commonly calls "antispam": filters, rules, to prune the set of results and filter out the majority that is of low quality and keep the minority of higher quality

- Anchor analysis

- Other graph-based approaches: TrustRank, etc. ...

# Context: Web Search Personalization

1. Ideas coming from **Recommender Systems**:
   - **Tailoring** the system's output to the current user.
   - Making "**recommendations**" of certain items vs. others to the user.
2. **Set of items =** all items returned by the IR engine
3. **Recommended items =** items that should be ranked higher / rank-merging with the IR score
4. Can be seen as another type of "**filter**"
5. **Content-based filtering:** based on user profile + item profile
6. **Collaborative filtering:** based on collaboratively collected info

**Link/Graph Analysis**          **vs**          **Recommender Systems**

# My Proposal: Previous Works Basis

- "A Large Scale Evaluation of Personalized Search Strategies", Dou, 2007
- "Topic-Sensitive PageRank", Haveliwala, 2002



8

# Project Definition

Features we want to achieve:

- **Web search personalization**: Results depending on the user

- **Collaborative approach:** results depending on users similar to current user

- **Usage-based approach:** the system adapts to what you do

- **Implicit and automated:** As a user, you do not need to do anything for the personalization to take place

**Data we have:** AOL Search Logs (2006), Internet-connection-reachable data

# Project Schedule(s)

## Original

| March | April | May | June | July | August | September |
|---|---|---|---|---|---|---|
| Subject definition | Literature review (2) | Thesis writing | | | Thesis written finalizatio | Finalization |
| Literature review (1) | System design | | | Evaluation | | Defense |
| | Minor development | Implementation | | | | |

## Effective

| March | April | May | June | July | August | September |
|---|---|---|---|---|---|---|
| Subject definition | Literature review (2) | | Thesis writing | | Thesis written finalization | |
| Literature review (1) | System design | | | | | Finalization |
| | Minor development | Implementation: System & Evaluation system | | | Evaluation (User Study) | Defense |

# Project Main Tasks

- User Model Definition

- Usage Extraction

- Collaboration

- Web Graph Personalized Scoring

- PageRank Personalization

- SERP Personalization

- Queries clustering $\quad c_i(q) = \dfrac{|kw(q, cluster(i))|}{|q|}$

$$c(q) = \begin{pmatrix} c_0(q) \\ c_1(q) \\ \vdots \\ c_{130}(q) \end{pmatrix}$$

- User profile $\quad c_l(u) = \displaystyle\sum_{p \in \mathcal{Q}(u)} P(q|u)w(q)c(q) \qquad P(q|u) = \dfrac{clicks(q, u)}{clicks(\bullet, u)}$

- Implementation: Runs fast enough. Using Numpy vectors.

# Personalization

- User-to-user similarity $\quad sim(u_1, u_2) = \dfrac{c_l(u_1) \cdot c_l(u_2)}{||c_l(u_1)||\ ||c_l(u_2)||}$

- Top 100 similar users (excluding sim = 1.0 ones)

- Scoring using similar users (collaboration) : $\quad score(u, q, p) = \dfrac{\sum\limits_{u_s \in \mathcal{S}_u} \left( sim(u_s, u) \cdot |clicks(q, p, u_s)| \right)}{\beta + \sum\limits_{u_s \in \mathcal{S}_u} |clicks(q, \bullet, u_s)|}$

- 3 implementations of the scoring:
  - **Straightforward**: store in DB, retrieve from DB when needed, with caching → Hugely² slow
  - **DB accesses grouped**: download & DB accesses by batches + caching → Hugely slow
  - **Store in RAM**, process in-RAM (no cache needed) → Quite OK
    - Multiprocessed in-RAM computation → Viable solution (~ 1 day) (10e-7s/sim)
    - Could be scaled using more CPUs / servers

13

# PageRank Personalization

- Recall: Standard PageRank $\rightarrow$ $R = c(M + E)R$
- Our *personalization vector*:

$$E(q, u) = \begin{pmatrix} p_0 \\ \vdots \\ p_i \\ \vdots \\ p_n \end{pmatrix}, p_i = \begin{cases} \dfrac{1}{N} & i \notin clicks(\mathcal{S}_u, \bullet, \bullet) \\ score(u, q, i) & i \in clicks(\mathcal{S}_u, \bullet, \bullet) \end{cases}$$

- CPPR formula $\rightarrow$ $CPPR(q, u) = c(M + E(q, u))CPPR(q, u)$

# AOL Re-Query / Web Crawl / Indexation

- AOL Re-Querying System
  - Loads keywords & related logs entries
  - Loads the SERP
  - Analyses SERP vs. logs to decide if we keep this SERP
  - Anti-bots protections workarounds: proxies, delays, tor, etc. ...

- Web Crawl
  - 7/3 other domains/same domain links following strategy

- Web Crawl Indexation
  - ElasticSearch with BM25
  - Several processing servers committing → central ElasticSearch Server

# User Study

- 5 queries

- 5 contexts (user + history)

- 11 volunteers

- Asked which preferred ranking

- Asked to select "at most 5 relevant links" for every ranking

- *Precision* metric: 
$$p_a(q, u) = \frac{\sum\limits_{r \in \mathcal{R}(q,u)} (11 - rank(r))}{\sum\limits_{i \in [[1,5]]} (11 - i)}$$
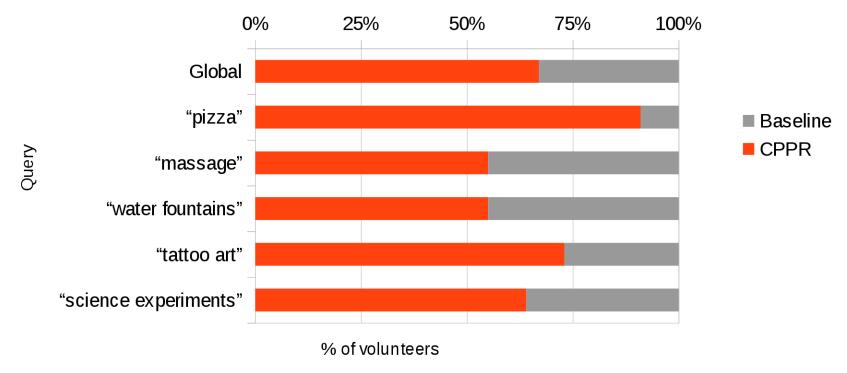
- Preferred ranking: CPPR

**NOT BAD**

"Best ranking overall" selection



Legend: Baseline (gray), CPPR (orange)

# Results (2)

- Average precision gain: 21.6%

### Over 11 Volunteers' Ratings

**>760**
commits

**1**
student

**200**GB
inter-servers exchanged data

**400**GB
cloud SSD peak usage

**45k**
written lines of codes

**95**
days github streak!

**700**
hours of work

**22**
CPUs peak usage

**1**
computer uptime record

**100**GB
laptop-uploaded data

**>300**GB
crawled websites

**11**
volunteers

**2M**
indexed documents

some figures might have been ceiled :)

**2**
cats pictures

# Technologies Used (that I already knew)

- Docker! (LXC) (now part of the "developer survival kit")

- CentOS, Fedora, Ubuntu, Debian

- Data Crunching: Python (http://python.org)

- Web Crawl: Scrapy (http://scrapy.org)

# Technologies Used (never used before) & learnt

- **Written work:** Multimarkdown
- **User Study Online Platform:** Express.js (http://expressjs.com)
- **PageRank computation:** Graph-Tool (http://graph-tool.skewed.de/)
- **HTML Parsing / Web Crawl Post-Processing:**
  - Chardet (https://github.com/chardet/chardet)
  - BeautifoulSoup4 (http://www.crummy.com/software/BeautifulSoup/)
- **IR/BM25 indexation & search:** ElasticSearch (http://elasticsearch.org)
- **Database:** MongoDB (http://mongodb.org)
- **Heavy computation:** Google Compute Engine (http://cloud.google.com/products/compute-engine/)
- **Python modules:** MultiProcessing, GZip, Pickle, JSON

# References / Bibliography

Please see the bibliography of the written thesis.

# THE END

Thanks for your attention.
Any questions?

# Web Crawl & Indexation (optional slide, for questions)

- Web Crawler based on Scrapy framework
- Follows links:
  - 3 links to same domain
  - 7 links to different domains
  - < 255chars
  - Some patterns excluded
  - Pictures, css, js, etc., excluded
- Trials with several different settings
- ~1M docs in <1d with "nice" settings (not hitting server too heavily)
- Indexation is another story: ~1.6 page/sec
  - Need for a several servers to get it to one day
  - Indexation process not built to be ran in parallel: most processes doing work that has already been done…
  - Separating data chunks by hand in the end...

# AOL Re-Query (optional slide, for questions)

- Web scraper
  - Loads keywords & related logs entries
  - Loads the SERP
  - Analyses SERP vs. logs to decide if we keep this SERP

- Anti-bots protections workarounds:
  - slow down → too slow, or banned
  - proxies → all banned
  - tor → banned
  - tor + slow down → too slow, or banned
  - proxies + slow down → tricky to add to the framework, but works OK