# CWI @ TREC 2013: Federated Web Search Track

**Alejandro Bellogín, Jiyin He, Arjen P. de Vries**

Centrum Wiskunde & Informatica, The Netherlands

{a.bellogin, j.he, arjen.de.vries}@cwi.nl

Text REtrieval Conference (TREC)
...to encourage research in information retrieval from large text collections.

## Resource selection

### ODP

Similarity between ODP's query and resource categories

Jaccard
Cosine

Search: www.arxiv.org
Open Directory Categories (1-10 of 10)
1. Science: Physics: Quantum Mechanics: Quantum Fi
2. Computers: Software: Operating Systems: Unix: BS
3. Computers: Internet: E-mail: Spam: Preventing (1)
4. Science: Math: Differential Equations: Dynamical S
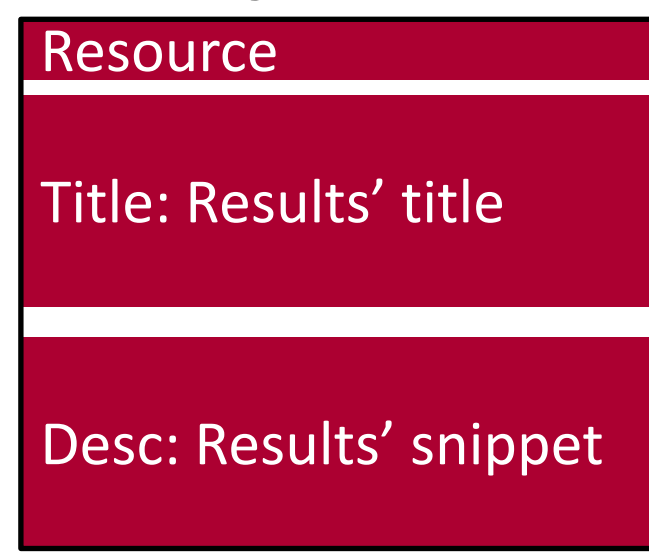5. Science: Math: Geometry: Computational Geometry

Search: retrieval
Open Directory Categories (1-25 of 100)
1. Computers: Software: Information Retrieval: Full
2. Computers: Software: Information Retrieval: Rank
3. Computers: Software: Information Retrieval (1B)
4. Reference: Knowledge Management: Knowledge R
5. Reference: Knowledge Management: Knowledge R

no order
no query text

### Retrieval model

Build pseudo-document and retrieve best matching resources

Resource

Title: Results' title

Desc: Results' snippet

Lucene

TF-IDF
BM25
Language Models

### Hybrid run

Aggregates rankings from the other methods using Borda voting

### FedWeb 2012

| Method | MAP | nDCG | MRR |
|---|---|---|---|
| TF-IDF+ODP Jacc | **0.338** | **0.516** | 0.564 |
| TF-IDF | 0.285 | 0.412 | **0.610** |
| ODP Jaccard | 0.283 | 0.471 | 0.439 |
| BM25 (1.2, 0.2) | 0.283 | 0.400 | 0.545 |
| LM ($\lambda = 0.1$) | 0.280 | 0.407 | 0.590 |
| ODP Cosine | 0.278 | 0.462 | 0.400 |
| BM25 (1.2, 0.8) | 0.272 | 0.397 | 0.557 |
| LM ($\lambda = 0.5$) | 0.263 | 0.394 | 0.571 |
| LM ($\lambda = 0.9$) | 0.252 | 0.387 | 0.566 |
| LM ($\lambda = 0.1$) desc | 0.241 | 0.386 | 0.602 |
| LM ($\mu = 200$) | 0.240 | 0.378 | 0.551 |
| LM ($\mu = 2000$) | 0.240 | 0.378 | 0.551 |
| BM25 (1.2, 0.8) desc | 0.239 | 0.383 | 0.608 |
| TF-IDF title | 0.215 | 0.321 | 0.495 |

### FedWeb 2013

Submitted best methods on 2012 collection

Results do not agree with 2013 collection

| Method | Run | nDCG@20 | ERR@20 |
|---|---|---|---|
| BM25 (1.2, 0.8) desc | - | **0.1588** | **0.0204** |
| LM ($\lambda = 0.1$) desc | - | 0.1476 | **0.0204** |
| BM25 (1.2, 0.2) | - | 0.1346 | 0.0068 |
| LM ($\lambda = 0.1$) | - | 0.1322 | 0.0068 |
| TF-IDF | CWI13SniTI | 0.1235 | 0.0067 |
| BM25 (1.2, 0.8) | - | 0.1223 | 0.0102 |
| LM ($\lambda = 0.5$) | - | 0.1218 | 0.0051 |
| LM ($\lambda = 0.9$) | - | 0.1153 | 0.0041 |
| LM ($\mu = 2000$) | - | 0.1033 | 0.0051 |
| TF-IDF title | - | 0.1016 | 0.0017 |
| TF-IDF+ODP Jacc | CWI13ODPTI | 0.0961 | 0.0034 |
| LM ($\lambda = 0.9$) | - | 0.0934 | 0.0017 |
| ODP Jaccard | CWI13ODPJac | 0.0497 | 0.0000 |

Best: Hybrid, TF-IDF, and ODP with Jaccard
Expected similar performance in 2013…

### FedWeb 2012

108 resources

Top 10 results
(snippets + pages)

Dec 2011 – Jan 2012

Queries

### FedWeb 2013

156 + 1 resources
(specific + BigWeb)

Top 10 results
(snippets + pages)

Apr – May 2013

Queries

| Best 2013 resource selection results | Method | P@10 | nDCG@20 | nDCG@50 | nDCG |
|---|---|---|---|---|---|
| | **2013 data** | | | | |
| | CWI13bstBM25desc* | **0.3408** | 0.1224 | 0.2024 | 0.5366 |
| | CWI13IndriQL | 0.3220 | **0.1622** | **0.2371** | **0.5438** |
| Best 2012 resource selection results | CWI13iaTODPJ | 0.2840 | 0.1509 | 0.1915 | 0.5253 |
| | CWI13bstTODPJ | 0.2500 | 0.1466 | 0.1839 | 0.4973 |
| | CWI13clTODPJ* | 0.1940 | 0.0551 | 0.0892 | 0.4610 |
| | **2012 data** | | | | |
| Also 2012! | CWI12bstTODPJ* | **0.4960** | 0.1246 | 0.1989 | 0.6081 |
| | CWI12IndriQL* | 0.4900 | **0.1464** | **0.2627** | **0.6525** |
| | CWI12clTODPJ* | 0.2200 | 0.0666 | 0.1106 | 0.5462 |
| | CWI12iaTODPJ* | 0.1940 | 0.0532 | 0.1015 | 0.5407 |

## Results merging

### Relevance

Documents ranked with respect to the query likelihood model:

$$p(d|q) \propto \prod_{w \in q} p(w|d)$$

Run: IndriQL

### Cluster

1. Rank resources (previous task)
2. Documents within a resource are ranked with IndriQL

Run: clTODPJ. Not submitted

### Diversity

IA-select diversification of IndriQL ranking using query relevance with respect to the resources

$$P(S|q) = \sum_z P(z|q)(1 - \prod_{d \in S}(1 - V(d|q, z)))$$

resource

Run: iaTODPJ

### Boost

Use directly the relevance with respect to the resources to boost the documents

$$p(d|q, z) \propto p(d|q)p(q|z)$$

Run: bstTODPJ

## Discussion

**Results merging can be solved with simple IR techniques**

Query likelihood obtains very good results

**How to define a training set for an evolving test environment?**

The rankings of the resources change

The content of the websites change

The type of queries is important:

are they tailored to be answered by a specific resource?

Distribution of ranking differences per resource