# Notes on Topic-sensitive PageRank

R. Basili

a.a. 2013-14

# The PageRank main workflow
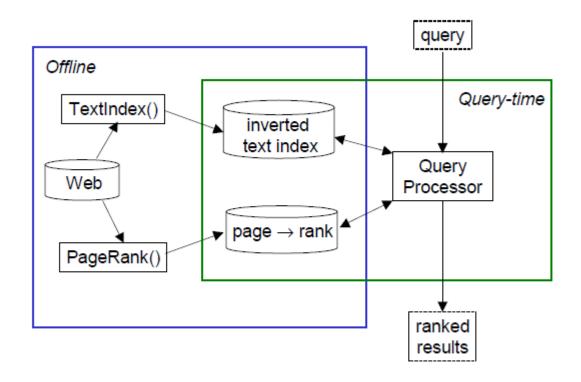


**Fig. 1.** Simplified diagram illustrating a simple search engine utilizing the standard PageRank scheme.

# Personalizing PageRank

- IDEA: determine the topics <span style="color:red">most closely associated with the query</span>, and use the appropriate topic-sensitive PageRank vectors for ranking the documents satisfying the query.

- This ensures that the "importance" scores reflect a preference for the link structure of pages that have some bearing on the query.

- As with ordinary PageRank, the topic-sensitive PageRank score can be used as part of a scoring function that takes into account other IR-based scores.

# Original PageRank

$$Rank = M' \times Rank$$
$$= (1 - \alpha)(M + D) \times Rank + \alpha p$$

# Topic-sensitive PageRank

Let $T_j$ be the set of URLs in the ODP category $c_j$. Then when computing the PageRank vector for topic $c_j$, in place of the uniform damping vector $p = [\frac{1}{n}]_{n \times 1}$, we use the nonuniform vector $p = v_j$ where

$$v_{ji} = \begin{cases} \frac{1}{|T_j|} & i \in T_j, \\ 0 & i \notin T_j. \end{cases} \qquad (7)$$

The PageRank vector for topic $c_j$ is given by $\boldsymbol{PR}(\alpha, \boldsymbol{v_j})$. We also generate the single unbiased PageRank vector (denoted as NoBias) for the purpose of comparison. The choice of $\alpha$ will

# Query Treatment

The second step in our approach is performed at query time. Given a query $q$, let $q'$ be the context of $q$. In other words, if the query was issued by highlighting the term $q$ in some Web page $u$, then $q'$ consists of the terms in $u$. Alternatively, we could use only those terms in $u$ nearby the highlighted term, as often times a single Web page may discuss a variety of topics. For ordinary queries not done in context, let $q' = q$. Using a multinomial naive-Bayes classifier [24],[6] with parameters set to their maximum-likelihood estimates, we compute the class probabilities for each of the 16 top-level ODP classes, conditioned on $q'$. Let $q'_i$ be the $i$th term in the query (or query context) $q'$. Then given the query $q$, we compute for each $c_j$ the following:

$$P(c_j|q') = \frac{P(c_j) \cdot P(q'|c_j)}{P(q')} \propto P(c_j) \cdot \prod_i P(q'_i|c_j) \tag{8}$$

# Combining Topics

Using a text index, we retrieve URLs for all documents containing the *original* query terms $q$. Finally, we compute the query-sensitive importance score of each of these retrieved URLs as follows. Let $r_{jd}$ be the rank of document $d$ given by the rank vector $PR(\alpha, v_j)$ (i.e., the rank vector for topic $c_j$). For the Web document $d$, we compute the query-sensitive importance score $s_{qd}$ as follows.

$$s_{qd} = \sum_j P(c_j|q') \cdot r_{jd} \qquad (9)$$

The results are ranked according to this composite score $s_{qd}$.[7]

The above query-sensitive PageRank computation has the following probabilistic interpretation, in terms of the "random surfer" model [26]. Let $w_j$ be the coefficient used to weight the $j$th rank vector, with $\sum_j w_j = 1$ (e.g., let $w_j = P(c_j|q)$). Then note that the equality

$$\sum_j [w_j PR(\alpha, v_j)] = PR(\alpha, \sum_j [w_j v_j]) \qquad (10)$$

holds, as shown in Appendix A. Thus we see that the following random walk on the Web yields the topic-sensitive score $s_{qd}$. With probability $1 - \alpha$, a random surfer on page $u$ follows an outlink of $u$ (where the particular outlink is chosen uniformly at random). With probability $\alpha P(c_j|q')$, the surfer instead jumps to one of the pages in $T_j$ (where the particular page in $T_j$ is chosen uniformly at random). The long term visit probability that the surfer is at page $v$ is exactly given by the composite score $s_{qd}$ defined above. Thus, topics exert influence over the final score in proportion to their affinity with the query (or query context).
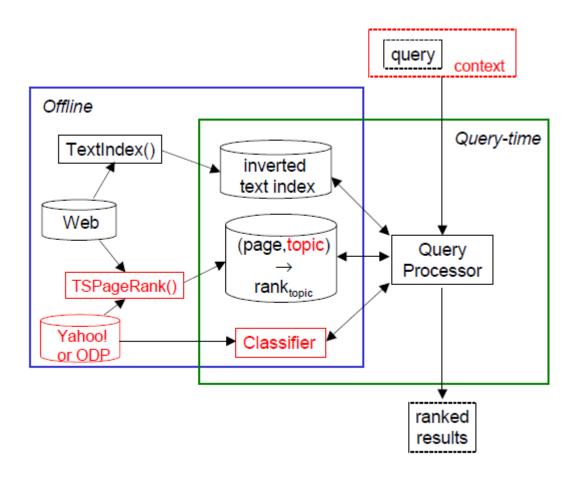
**Fig. 2.** Illustration of our system utilizing topic-sensitive PageRank.

**Table 3.** Topic pairs yielding most similar rankings.

| Bias-Topic Pair | $OSim$ | $KSim$ |
|---|---|---|
| (GAMES, SPORTS) | 0.18 | 0.13 |
| (NOBIAS, REGIONAL) | 0.18 | 0.12 |
| (KIDS & TEENS, SOCIETY) | 0.18 | 0.11 |
| (HEALTH, HOME) | 0.17 | 0.12 |
| (HEALTH, KIDS & TEENS) | 0.17 | 0.11 |

**Table 5.** Estimates for $P(c_j|q)$ for a subset of the test queries.

| alcoholism | |
|---|---|
| HEALTH | 0.47 |
| KIDS & TEENS | 0.20 |
| ARTS | 0.06 |

| bicycling | |
|---|---|
| SPORTS | 0.52 |
| REGIONAL | 0.13 |
| HEALTH | 0.07 |

| blues | |
|---|---|
| ARTS | 0.52 |
| SHOPPING | 0.12 |
| NEWS | 0.08 |

| citrus groves | |
|---|---|
| SHOPPING | 0.34 |
| HOME | 0.21 |
| REGIONAL | 0.18 |

| classical guitar | |
|---|---|
| ARTS | 0.75 |
| SHOPPING | 0.21 |
| NEWS | 0.01 |

| computer vision | |
|---|---|
| COMPUTERS | 0.24 |
| BUSINESS | 0.14 |
| REFERENCE | 0.09 |

| cruises | |
|---|---|
| RECREATION | 0.65 |
| REGIONAL | 0.18 |
| SPORTS | 0.04 |

| death valley | |
|---|---|
| REGIONAL | 0.28 |
| SOCIETY | 0.14 |
| NEWS | 0.10 |

| field hockey | |
|---|---|
| SPORTS | 0.89 |
| SHOPPING | 0.03 |
| REFERENCE | 0.03 |

| graphic design | |
|---|---|
| COMPUTERS | 0.36 |
| BUSINESS | 0.23 |
| SHOPPING | 0.09 |

| gulf war | |
|---|---|
| SOCIETY | 0.21 |
| KIDS & TEENS | 0.18 |
| REGIONAL | 0.17 |

| hiv | |
|---|---|
| HEALTH | 0.40 |
| NEWS | 0.19 |
| KIDS & TEENS | 0.14 |

| java | |
|---|---|
| COMPUTERS | 0.53 |
| GAMES | 0.10 |
| KIDS & TEENS | 0.06 |

| lyme disease | |
|---|---|
| HEALTH | 0.96 |
| REGIONAL | 0.01 |
| RECREATION | 0.01 |

| mutual funds | |
|---|---|
| BUSINESS | 0.77 |
| REGIONAL | 0.05 |
| HOME | 0.05 |

| parallel architecture | |
|---|---|
| COMPUTERS | 0.70 |
| SCIENCE | 0.10 |
| REFERENCE | 0.07 |

| rock climbing | |
|---|---|
| RECREATION | 0.54 |
| REGIONAL | 0.13 |
| SPORTS | 0.07 |

| san francisco | |
|---|---|
| SPORTS | 0.27 |
| REGIONAL | 0.16 |
| RECREATION | 0.10 |

| shakespeare | |
|---|---|
| ARTS | 0.34 |
| REFERENCE | 0.21 |
| KIDS & TEENS | 0.15 |

| table tennis | |
|---|---|
| SPORTS | 0.53 |
| SHOPPING | 0.14 |
| REGIONAL | 0.09 |

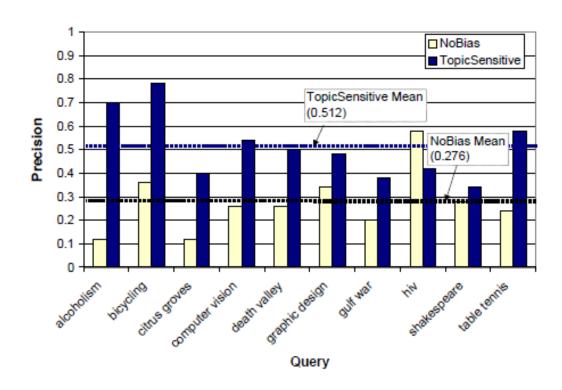| telecommuting | |
|---|---|
| BUSINESS | 0.70 |
| KIDS & TEENS | 0.04 |
| SOCIETY | 0.03 |

**Fig. 3.** Precision @ 10 results for our test queries. The average precision over the ten queries is also shown.

# Applications

- Domain specific PageRank
- User Profiling, when **p** is based on user preferred topics
- Semantic Classification:
  - Word Sense Disambiguation (Soroa & Agirre, 2009) when:
    - The graph is a sense dictionary (i.e. Wordnet)
    - **PR** gives the preference to a sense as its reachability
    - **p** is based on the context of the target word

# References

- Taher H. Haveliwala, *Topic-Sensitive PageRank* In Proceedings of the Eleventh International World Wide Web Conference, May 2002.

- E. Agirre and A. Soroa. 2009. *Personalizing pagerank for word sense disambiguation*. in Proceedings of EACL-09, Athens, Greece.