

# TOPIC-SENSITIVE PAGE RANK

- choose independent categories in advance

$C_1, C_2, \dots, C_k$

example: } business  
                  } economics  
                  } sports  
                  } news  
                  } travel...

- compute page rank for each page  $u$  inside a set of "relevant" pages for each category  $k$

$$PR_k(u) = PR \left\{ \begin{array}{l} \text{of page } u \text{ in set} \\ \text{retrieved for category } k \end{array} \right\}$$

- at query time  $q$

$$PR_q(u) = \sum_k \text{prob}(C_k | q) \cdot PR_k(u)$$

-  $\text{prob}(C_k | q) \approx \underbrace{\text{prob}(C_k)}_{\substack{\text{estimate} \\ \text{in} \\ \text{advance}}} \cdot \underbrace{\text{prob}(q | C_k)}_{\text{language model.}}$

The *outlinking* of a node  $v$  is the set of all links selected by the predicate  $P$  from  $v$  to  $u$ :

$$O_v^P = \{ (v, u) \in E : P(v, u) \}$$

The *inlinking* of a node  $u$  is the set of all links selected by the predicate  $P$  from  $v$  to  $u$ :

$$I_u^P = \{ (v, u) \in E : P(v, u) \}$$

The *base set*  $B^P$  of a node  $v$  is the set  $B^P$  and  $A^P$  is defined to be:

$$B_v^P = \{ u \in V : (v, u) \in O_v^P \}$$

The *neighborhood*  $N_v^P$  of a node  $v$  is the base set  $B_v^P$  as its vertex set and an edge set  $N_v^P$  containing those edges in  $E$  that are covered by  $P$  and pointed by  $v$ :

$$N_v^P = \{ (u, v) \in E : (v, u) \in O_v^P \wedge P(v, u) \}$$

To simplify notation, we write  $B$  to denote  $B_v^P$  and  $N$  to denote  $N_v^P$ .

### 5.2 The HITS algorithm

For each node  $v$  in the neighborhood graph, HITS computes two scores: a *hub score*  $H(v)$ , estimating how authoritative a node  $v$  is, and a *authority score*  $A(v)$ , estimating how good a node  $v$  is to many authoritative pages. This is done using the following algorithm:

**HITS-1Hub-and-Authority-Scores:**

1. For all  $u \in B$  do  $H(u) := \sqrt{\|A\|_2} \cdot \|u\|_2 := \sqrt{\frac{1}{|B|}}$ .
2. Repeat until  $H$  and  $A$  converge:
  - (a) For all  $u \in B : A'(u) := \sum_{(v,u) \in N} H(v)$
  - (b) For all  $u \in B : H'(u) := \sum_{(u,v) \in N} A(v)$
  - (c)  $H := \frac{1}{\|A'\|_2} H', A := \frac{1}{\|H'\|_2} A'$

where  $\|X\|_2$  is the euclidean norm of vector  $X$ .

### 5.3 The SALSA algorithm

For each node  $u$  in the neighborhood graph, SALSA computes an *authority score*  $A(u)$  and a *hub score*  $H(u)$  using the following two independent algorithms:

**SALSA-Hub-Scores:**

1. Let  $B^H$  be  $\{ u \in B : out(u) > 0 \}$ .
2. For all  $u \in B$

$$H(u) = \begin{cases} \frac{1}{|B^H|} & u \in B^H \\ 0 & \text{otherwise} \end{cases}$$

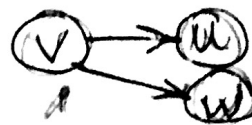
3. Repeat until  $H$  converges:
  - (a) For all  $u \in B^H$ :

$$H'(u) = \sum_{(u,v) \in N} \frac{H(v)}{in(v) \cdot out(v)}$$

- (b) For all  $u \in B^H : H(u) = H'(u)$



$$H(u) = \sum_v \sum_w \frac{H(w)}{out(w) \cdot in(v)}$$



$$A(w) = \sum_{v \in V} \frac{A(v)}{in(w) \cdot out(v)}$$

$$B^A = \{ u \in B : in(u) > 0 \}$$

$$A(u) = \begin{cases} \frac{1}{|B^A|} & u \in B^A \\ 0 & \text{otherwise} \end{cases}$$

Repeat until  $A$  converges:

- (a) For all  $u \in B^A$

$$A'(u) = \sum_{(v,u) \in N} \sum_{(v,w) \in N} \frac{A(w)}{out(v) \cdot in(w)}$$

- (b) For all  $u \in B^A : A(u) = A'(u)$

## 6. THE SCALABLE HITS LINK STORE

We have described the HITS algorithm on the *Scalable Neighborhood Graph* for the purpose of illustrating the design of the HITS Link Store. The HITS Link Store is a distributed system that stores the neighborhood graph in main memory and efficiently distributes access to nodes (URLs) and edges (links) and that uses data compression techniques that preserve structural properties (namely, the presence of related links) in the neighborhood graph to achieve fairly good compression. Storing the full 17.7 billion link graph in a networked system requires six machines, each with 16 GB of main memory.

The two principal abstractions used in HITS are a *URL store* and two *link stores*, one to trace links forward and another to trace them back. Clients use HITS by linking against a library containing classes that implement clerks for the URL store and the link stores; all the intricacies common to distributed systems are handled by the clerks and the HITS servers.

The URL store maintains a bijection between URLs (strings) and UID (integers) that serves as a handle for URLs. Clients can map URLs to UIDs and UIDs back to URLs. The API of the URL store clerk looks as follows:

```
class UrlStoreClerk {
    ... // omitting private members
public:
    UrlStoreClerk(char *serverNameFile);
    UrlStoreClerk();
    INT64 UrlToUid(char *url);
    char *UidToUrl(INT64 uid);
    SeqInt64 BatchedUrlToUid(SeqString& urls);
    SeqString BatchedUidToUrl(SeqInt64& uids);
    ... // omitting methods irrelevant to this paper
};
```

The *UrlStoreClerk* constructor takes the name of a file that contains the names of the HITS servers maintaining the graph. The central methods are *UrlToUid*, which maps a URL to a UID, and *UidToUrl*, which maps a UID back to a URL. The methods *BatchedUrlToUid* and *BatchedUidToUrl* are variants of the previous two methods that allow the mapping of entire batches of URLs or UIDs; their purpose is to allow client applications to amortize RPC overheads. As a point of reference, mapping a URL to a UID takes about 3 microseconds, while performing a null RPC takes about 100 microseconds; so providing a mechanism to batch up requests is performance-critical. Our implementations of