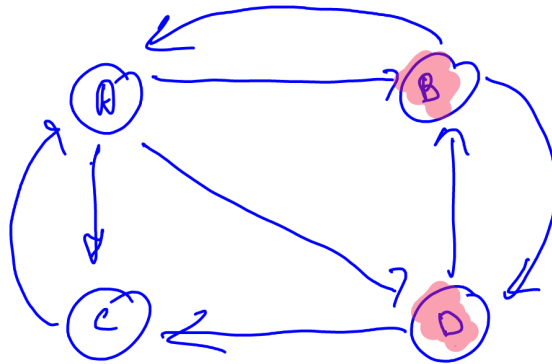# Topic - Sensitive   Page Rank



$$M = \begin{bmatrix} 0 & 1/2 & 1 & 0 \\ 1/3 & 0 & 0 & 1/2 \\ 1/3 & 0 & 0 & 1/2 \\ 1/3 & 1/2 & 0 & 0 \end{bmatrix}$$

$\beta = 0.8 \ / \ 0.85 \ / \ 0.3$

$$v' = \beta M v + (1-\beta) \boxed{e}$$

$$e = \begin{pmatrix} 1/n \\ \vdots \\ 1/n \end{pmatrix} \longrightarrow e_S = \frac{1}{|S|} \begin{pmatrix} 0 \\ 1 \\ 0 \\ 1 \\ 0 \end{pmatrix} \begin{matrix} \\ \leftarrow \\ \\ \leftarrow \\ \end{matrix} \ \ --- \ \underset{B \ and \ P}{\text{sports}}$$

$$S = \{ B, D \} \longrightarrow e_S = \frac{1}{2} \begin{pmatrix} 0 \\ 1 \\ 0 \\ 1 \end{pmatrix}$$

$$v' = \begin{bmatrix} 0 & 2/5 & 4/5 & 0 \\ 4/15 & 0 & 0 & 2/5 \\ 4/15 & 0 & 0 & 2/5 \\ 4/15 & 2/5 & 0 & 0 \end{bmatrix} v + \begin{bmatrix} 0 \\ 1/10 \\ 0 \\ 1/10 \end{bmatrix} \qquad \beta = 0.8$$

$$e_S = \begin{bmatrix} 0 \\ ^1/2 \\ 0 \\ ^1/2 \end{bmatrix} , \begin{bmatrix} ^2/20 \\ ^3/20 \\ ^2/10 \\ ^3/10 \end{bmatrix} \cdots \begin{bmatrix} ^{54}/210 \\ ^{59}/210 \\ ^{38}/210 \\ ^{59}/290 \end{bmatrix}$$

"Usual PageRank + Taxation"

$$\begin{bmatrix} ^{15}/998 \\ ^{19}/998 \\ ^{35}/998 \\ ^{19}/998 \end{bmatrix}$$

## Importance of Words in Documents

Topics $\longrightarrow$ special words

baseball $\rightarrow$ "ball", "hit", "run", "pitch"

Find significant words?

$\rightarrow$ most frequent words?    stop words ("the", "and"...)

Indicators of topics = relatively rare words

"notwithstanding" / "albeit"

Polo $\rightarrow$ "chakker"

TF · IDF ( Term Frequency times Inverse Document Frequency)

$$TF_{ij} = \frac{f_{ij}}{\max_K f_{Kj}} \quad \text{(frequency word i at doc j, exp. stop-words)}$$

word $\rightarrow ij \leftarrow$ document

$N$ documents in collection

term $i$ appears in $n_i$ documents

$$\boxed{IDF_i = \log_2 \left( N/n_i \right)}$$

For term $i$: $\boxed{TF_{ij} \times IDF_i}$
and doc $j$

Example    $20^{20}$ documents, word $w$ appears in $2^{10}$

$$IDF_w = \log_2 \left( 2^{20}/2^{10} \right) = \log_2 \left( 2^{10} \right) = 10.$$

If $w$ appears in the doc $j$ 20 times

$$TF_{wj} = 1 \quad \rightarrow \quad TF \cdot IDF_{wj} = 10$$

document $k$, $w$ appears only 1

$$TF_{wk} = \frac{1}{20} \quad \rightarrow \quad TF \cdot IDF_{wj} = \frac{10}{20} = \frac{1}{2},$$

---

If the word "$f$" appears in each doc.
what will be $IDF_f$ ? $IDF_f = 0$