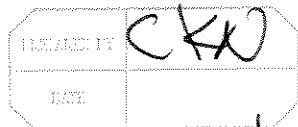2/8/07

B-Y & R-N 2.5.4

Notes on Probabilistic IR

P. 1/

Given a query $q$, and some documents $d_1, ..., d_n$, we want to estimate the probability

$$P(\ d_i \text{ is relevant to } q)$$

and perhaps rank the documents with respect to this probability.

How to do this? Assume the query and documents are made up of terms $t_i$

idea: terms that occur in known relevant docs more often than in the collection as a whole should indicate higher (probability of) relevance in documents whose relevance (yes, no) is still unknown.

user interaction may be used, but isn't necessary (!)

Assumption 1: probability of doc $d_i$ being relevant to query $q$ depends only on $d_i$ and $q$, not on other documents or queries or whatever

Assumption 2: a certain subset, $R$, of the collection exists which the user will consider relevant.

$$R = \{d_i, \ldots, d_k\}$$ ↑relevant docs
or we think so...

$\overline{R}$ is the rest of the corpus

$P(R|d_i)$ is the probability of $d_i$ being in $R$ (∴ relevant)

$P(\overline{R}|d_i)$ prob. of irrelevance

$$\text{sim}(d_i, q) = \frac{P(R|d_i)}{P(\overline{R}|d_i)}$$

unfortunately, we don't know exactly which documents are in $R$, nor even how many such documents there may be

i.e. $P(d_i$ is relevant$)$

$$P(R \mid d_i) = \frac{P(d_i \mid R)\, P(R)}{P(d_i)}$$

i.e. $P(d_i$ is relevant$)$

$$P(\bar{R} \mid d_i) = \frac{P(d_i \mid \bar{R})\, P(\bar{R})}{P(d_i)}$$

the denominators cancel out, and we have

$$sim(d_i, g) = \frac{P(d_i \mid R)\, P(R)}{P(d_i \mid \bar{R})\, P(\bar{R})}$$

$P(R)$ (and therefore $P(\bar{R})$) are the same for all documents, so for ranking of documents

$$sim(d_i, g) \sim \frac{P(d_i \mid R)}{P(d_i \mid \bar{R})}$$

Assumption 3: index terms are independent, so...

$sim(d_i, q) \sim$

$$\frac{\prod_{\substack{all\ k_j \\ in\ d_i}} P(k_j|R) \times \prod_{\substack{all\ k_j \\ not\ in\ d_i}} P(\overline{k_j}|R)}{\prod_{\substack{all\ k_j \\ in\ d_i}} P(k_j|\overline{R}) \times \prod_{\substack{all\ k_j \\ not\ in\ d_i}} P(\overline{k_j}|\overline{R})}$$

$k_j$ in random doc in R

$k_j$ not in $d_i$

$P(k_j|R) + P(\overline{k_j}|R) = 1$
since every keyword will be in
a random document in R, (or in $\overline{R}$)
it won't — for sure. So we
have

$sim(d_i, q) \sim$

$$\frac{\prod P(k_j|R) \times \prod (1 - P(k_j|R))}{\prod P(k_j|\overline{R}) \times \prod (1 - P(k_j|\overline{R}))}$$

taking logs, noting that
$\log \prod P(k_j|R) = \sum \log P(k_j|R)$

we have

after some algebra
(and magic)

$$\text{sim}(d_i, q) \sim$$

$$\sum w_{i,q} \cdot w_{i,j} \cdot \left( \log \frac{P(k_i \mid R)}{1 - P(k_i \mid R)} + \right.$$

$$\left. \log \frac{1 - P(k_i \mid \bar{R})}{P(k_i \mid \bar{R})} \right)$$

$P(k_i \mid R)$ can start at 0.5

$P(k_i \mid \bar{R})$ can start at

$n_i / N$, where

$n_i = \#$ of docs containing $k_i$

$N = \#$ of docs in collection