

Retrieval by Content

Part 3: Text Retrieval Latent Semantic Indexing

Latent Semantic Indexing (LSI)

- Disadvantage of exclusive use of representing a document as a T-dimensional vector of term weights
 - Users may pose queries using terms different from terms used to index a document
 - E.g., term *data mining* is semantically similar to *knowledge discovery*

LSI method

- Approximate the T-dimensional term space by k principal components directions in this space
 - Using the $N \times T$ document term matrix to estimate directions
 - Results in a $N \times k$ matrix
 - Terms *database*, *SQL*, *indexing* etc are combined into a single principal component

Singular Value Decomposition

- Find decomposition of $N \times T$ document-term matrix M as follows:
- $M=USV^T$

$N \times T$

$T \times T$
Diagonal Matrix
of eigen values of
principal directions

$T \times T$ matrix whose columns are new
Orthogonal bases for the data

Singular Value Decomposition

Document-Term Matrix, M

	database	SQL	index	regression	likelihood	linear
D1	24	21	9	0	0	3
D2	32	10	5	0	3	0
D3	12	16	5	0	0	0
D4	6	7	2	0	0	0
D5	43	31	20	0	3	0
D6	2	0	0	18	7	16
D7	0	0	1	32	12	0
D8	3	0	0	22	4	2
D9	1	0	0	34	27	25
D10	6	0	0	17	4	23

Find a decomposition $M = USV^T$

U is a 10 x 6 matrix of weights
(each row for a particular document)
S is a 6 x 6 diagonal matrix of Eigen values
Columns of 6 x 6 matrix V^T represent
principal components (or orthogonal bases)

S matrix has diagonal elements
77.4, 69.5, 22.9, 13.5, 12.1, 4.8

Variance is captured by first two
Elements. Fraction of variance captured is

$$\frac{\lambda_1^2 + \lambda_2^2}{\sum_i \lambda_i^2} = 0.925 \text{ or only } 7.5\% \text{ of data is lost}$$

U Matrix (using 2 PCs)

Document	PC1	PC2
d1	30.8998	-11.4912
d2	30.3131	-10.7801
d3	18.0007	-7.7138
d4	8.3765	-3.5611
d5	52.7057	-20.6051
d6	14.2118	21.8263
d7	10.8052	21.914
d8	11.508	28.0101
d9	9.5259	17.7666
d10	19.9219	45.0751

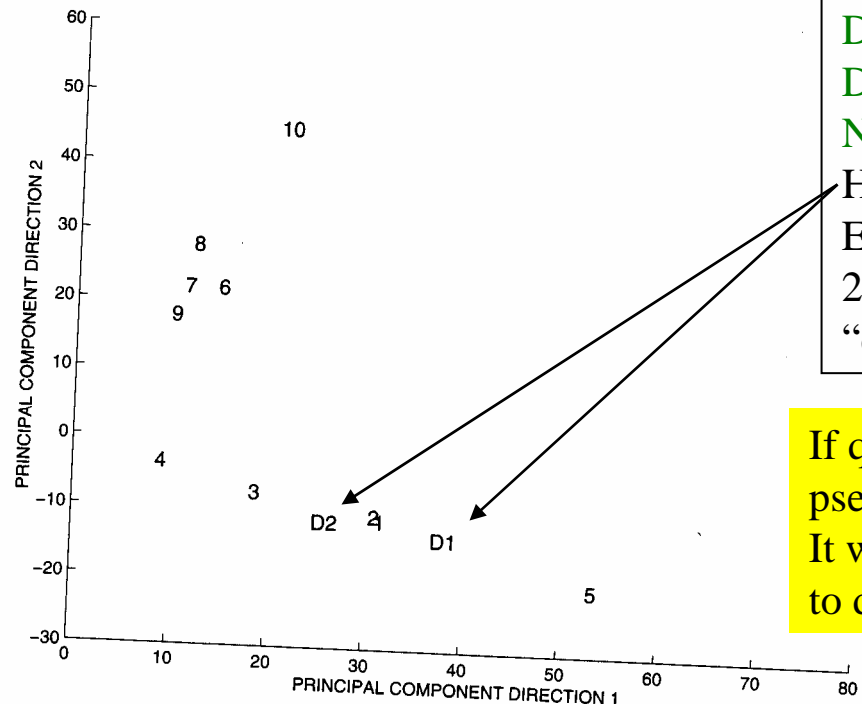
V matrix

	database	SQL	index	regression	likelihood	linear
v1	0.74	0.49	0.27	0.28	0.18	0.19
v2	-0.28	-0.24	-0.12	0.74	0.37	0.31

Two directions in which data is most spread out. First emphasizes *database* and *SQL*. Second emphasizes *regression*, *likelihood*, *linear*

LSI Method: First Two Principal Components of Document Term Matrix

Emphasizes
Regression, Likelihood, Linear



D1: has *database* 50 times
 D2: has *SQL* 50 times
 None of the other terms
 Have small distance in LSI
 Even though each is missing
 2 of 3 terms associated with
 “database” direction.

If query is *SQL*, with
 pseudoterm representation:
 It will be closer in angle
 to database direction

Emphasizes
Database, SQL

Figure 14.3 Projected locations of the 10 documents (from table 14.2) in the two dimensional plane spanned by the first two principal components of the document-term matrix M .

LSI Practical Issues

- Query is represented as a vector in PCA space and angle calculated
 - E.g., Query *SQL* is converted into pseudo vector
- In practice, computing PCA vectors directly is computationally infeasible.
 - Special purpose sparse SVD techniques for high- dimensions are used
- Can also model Document-Term matrix probabilistically as a mixture of simpler component distributions
 - Each component represents distribution of terms conditioned on a particular topic
 - Each component can be a naïve Bayes model

Incorporating User Feedback in Document Retrieval

- Retrieval Algorithms have a more interactive flavor than other data mining algorithms
- A user with query Q may be willing to iterate through a few sets of different retrieval trials and provide user feedback to the algorithm by labeling returned documents as relevant and non-relevant
- Applicable to any retrieval system not just text retrieval

Relevance Feedback

- Principle: Relevance is user centric
- If user could see all documents
 - user could separate them into two sets *relevant* R and *non-relevant* NR
 - Second round of input is called *Relevance Feedback*
- Goal is to learn from these sets to refine results
- Given these two sets, the optimal query is

$$Q_{optimal} = \frac{1}{|R|} \sum_{D \in R} D - \frac{1}{|NR|} \sum_{D \in NR} D$$

- Where D is a term-vector representation for documents

Rocchio's Algorithm

- Assume user has not used optimal query
- Instead has a specific query $Q_{current}$
- Algorithm uses this to return a small set of documents which are **labeled by user** as relevant R' and non-relevant NR'
- Rocchio's algorithm refines the query thus:

$$Q_{new} = \alpha Q_{current} + \frac{\beta}{|R'|} \sum_{D \in R'} D - \frac{\gamma}{|NR'|} \sum_{D \in NR'} D$$

where α , β and γ are heuristically chosen constants that control sensitivity to most recent labeling

Query is modified by

moving current query toward mean vector of documents judged relevant
and away from those considered irrelevant

Process is repeated with user again labeling documents

Pseudo Relevance Feedback

$$Q_{new} = \alpha Q_{current} + \frac{\beta}{|R'|} \sum_{D \in R'} D - \frac{\gamma}{|NR'|} \sum_{D \in NR'} D$$

- Collect R' assuming certain number of most relevant documents are relevant
 - γ is set to zero
 - Typically top 10 to 20 are used

Probabilistic Relevance Feedback

- Tune retrieval system to a statistical model of the generation of documents and queries
- Method of ranking documents is based on an odds ratio for relevance
- Let R be a Boolean value indicating relevance of document D wrt query q

$$\begin{aligned} & \frac{P(R/q, D)}{P(NR/q, D)} \\ &= \frac{P(R, q, D) / P(q, D)}{P(NR, q, D) / P(q, D)} \\ &= \frac{P(R/q)P(D/R, q)}{P(NR/q)P(D/NR, q)} \end{aligned}$$

Use Naïve Bayes model
where terms are assumed independent

Naïve Bayes model of Probabilistic Retrieval

$$\frac{P(D / R, q)}{P(D / NR, q)} = \prod_t \frac{P(x_i / R, q)}{P(x_i / NR, q)}$$

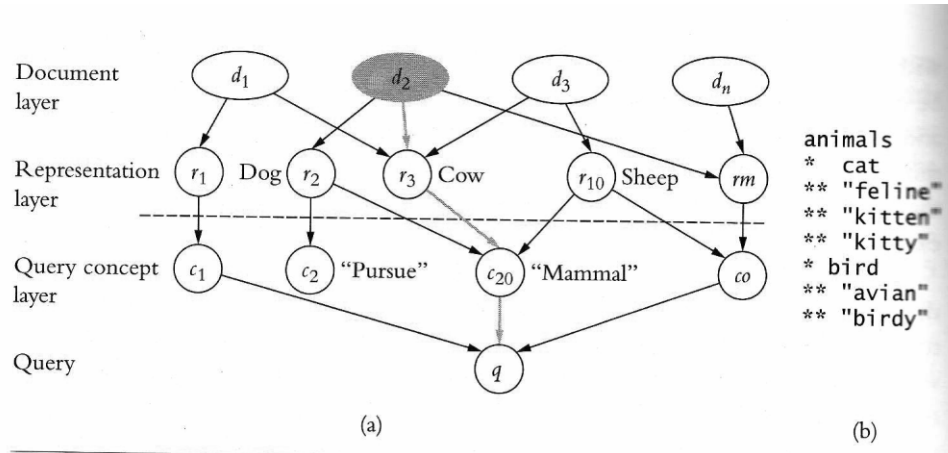
- Let $a_{t,q} = P(x_i=1/R, q)$ and $b_{t,q} = P(x_i=1/NR, q)$ since the terms are present/absent, i.e., the features are binary-valued
- Hence, the standard two-class independent binary classification result holds:

$$\frac{P(D / R, q)}{P(D / NR, q)} \propto \prod_t \frac{a_{t,q} (1 - b_{t,q})}{b_{t,q} (1 - a_{t,q})}$$

- Parameters $a_{t,q}$ and $b_{t,q}$ have to be estimated
- Disadvantage: User has to rate some responses before probabilities kick in

Other Probabilistic Models

- Bayesian inference network
 - Nodes correspond to documents, terms, “concepts” and queries



- Most IR systems in use today use standard vector-space models and not probabilistic retrieval models

Automated Recommender Systems

- Instead of modeling preferences of a single user, generalize to the case where there is information about multiple users
- Collaborative Filtering
 - Method to leverage group information
 - Example: you purchase a CD at a website
 - Algorithm provides list of CDs others who also purchased that CD
 - Generalize based on user profile: need
 - vector representation
 - Similarity metrics