Latent Semantic Indexing

Thanks to Ian Soboroff

Information Retrieval

1

Issues: Vector Space Model

 Assumes terms are independent • Some terms are likely to appear together synonyms, related words spelling mistakes? Terms can have different meanings depending on context Term-document matrix has very high dimensionality are there really that many important features for each document and term? Lecture 12 Information Retrieval 2





- Compute singular value decomposition of a term-document matrix
 - D, a representation of M in *r* dimensions

Lecture 12

- T, a matrix for transforming new documents
- diagonal matrix Σ gives relative importance of dimensions

Information Retrieval

3

LSI Term matrix T

T matrix

- gives a vector for each term combo in LSI space
- for a new document c, c'*T gives a new row in D
- That is, "fold in" the new document into the LSI space

• LSI is a rotation of the term-space

- original matrix: terms are d-dimensional
- new space has lower dimensionality
- dimensions are groups of terms that tend to cooccur in the same documents
 - synonyms, contextually-related words, variant endings

Lecture 12











Truncating Dimensions in LSI

	3-grams, tl.idl	5-grams, tl.idl	
0.45 0.4 - 0.35 - 0.3 - 0.3 - 0.25 - 0.2 - 0.15 - 0.1 -+		0.5 0.45 0.45 0.45 0.4 0.4 0.35 0.3	
0.05	5 10 15 20 25	0.05 <u>1 1 1 1</u> 0 5 10 15 20	25
0.24 0.22 0.22 0.2 0.2 0.18 0.16 0.16	3-grams, ti only	5-grams, 11 only 0.35 0.35 0.3 0.35 0.35 0.25 ++++++++++++++++++++++++++++++++++++	
0.12 0.1 0.08 0	5 10 15 20 25 Dimension	0.15 0.05 0.05 0 5 10 15 20 Dimension	
Lecture 12	Informatio	on Retrieval	9

Document matrix D

- D matrix
 - coordinates of documents in LSI space
 - same dimensionality as T vectors
 - can compute the similarity between a term and a document



Improved Retrieval with LSI

- New documents and queries are "folded in"
 - multiply vector by $T\Sigma^{-1}$
- Compute similarity for ranking as in VSM
 - compare queries and documents by dot-product
- Improvements come from
 - reduction of noise
 - no need to stem terms (variants will co-occur)
 - no need for stop list
 - stop words are used uniformly throughout collection, so they tend to appear in the first dimension
 - No speed or space gains, though...

Lecture 12

LSI in TREC-3

- LSI space computed from a sample of the document collection
- Documents and queries folded into LSI space for comparison
- Improvement in AP with LSI: 5%
 - Improvements up to 20% seen in smaller collections

Other LSI Applications

- Text classification
 - by topic
 - dimension reduction -> good for clustering
 - by language
 - languages have their own stop words
 - by writing style
- Information Filtering
- Cross-language retrieval

N-gram indexing recap

- Index all *n* character sequences
 - language-independent
 - resistant to noisy text
 - no stemming
 - easy to do
- Document ⇒ array of n-gram frequencies

n=5World Hello Hello World Hello World Hello World

14

Lecture 12

Information Retrieval

Why N-grams?

- N-grams capture pairs of words
 - Brings out phraseology and word choice
- LSI using n-grams might cluster documents by writing style and/or author
 - a lot of what makes style is word choices and stop word usage
- Small experiment
 - Three biblical Hebrew texts: Ecclesiastes, Song of Songs, Book of Daniel

15

• used 3-grams in original Hebrew







Conclusion

- LSI can be a useful technique for reducing the dimensionality of an IR problem
 - reduction can improve effectiveness
 - reduction can find surprising relationships!
- SVD can be expensive to compute on large matrices
- Available tools for working with LSI
 - MATLAB or Octave (small data sets only)
 - SMART (an IR system) with SVDPACK