

# Retrieval by Content

## Part 2: Text Retrieval

### Term Frequency and Inverse Document Frequency

# Text Retrieval

- Retrieval of text-based information is referred to as Information Retrieval (IR)
- Used by text search engines over the internet
- Text is composed of two fundamental units documents and terms
- Document: journal paper, book, e-mail messages, source code, web pages
- Term: word, word-pair, phrase within a document

# Representation of Text

- Capability to retain as much of the semantic content of the data as possible
- Computation of distance measures between queries and documents efficiently
- Natural language processing is difficult, e.g.,
  - *Polysemy* (same word with different meanings)
  - *Synonymy* (several different ways to describe the same thing)
- IR systems in use today do not rely on NLP techniques
  - Instead rely on vector of term occurrences

# Vector Space Representation

Document-Term Matrix

	t1	t2	t3	t4	t5	t6
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

Terms used  
In "Database"  
Related

t1	database
t2	SQL
t3	index
t4	regression
t5	likelihood
t6	linear

"Regression" Terms

$d_{ij}$  represents number of times  
that term appears in that document

# Cosine Distance between Document Vectors

Document-Term Matrix

	t1	t2	t3	t4	t5	t6
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

Cosine Distance

$$d_c(D_i, D_j) = \frac{\sum_{k=1}^T d_{ik} d_{jk}}{\sqrt{\sum_{k=1}^T d_{ik}^2 \sum_{k=1}^T d_{jk}^2}}$$

Cosine of the angle between two vectors

Equivalent to their inner product after each has been normalized to have unit

Higher values for more similar vectors

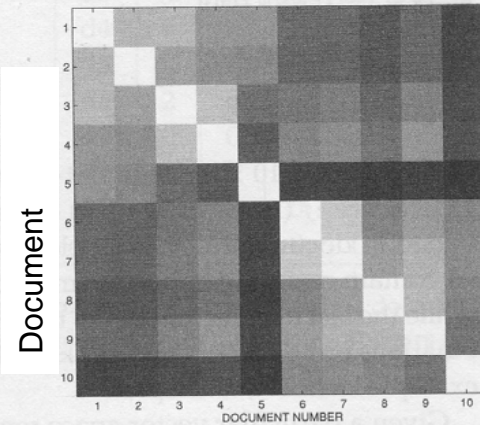
Reflects similarity in terms of relative distributions of components

Cosine is not influenced by one document being small compared to the other (as in the case of Euclidean)

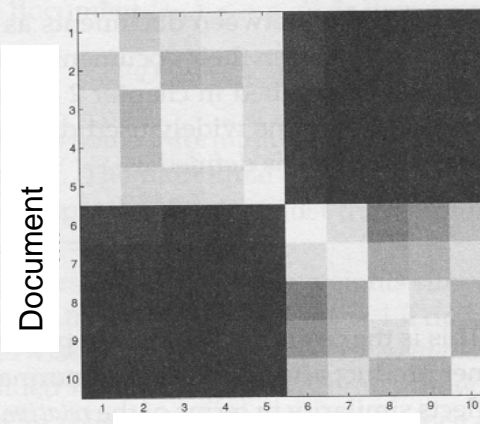
# Euclidean vs Cosine Distance

## Document-Term Matrix

	t1	t2	t3	t4	t5	t6	
Database	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
Regression	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



**Euclidean**  
White = 0  
Black = max distance



**Cosine**  
White = Larger Cosine (or smaller angle)

Both show two clusters of light sub-blocks (*database* documents and *regression* documents)

**Euclidean**: 3,4 closer to 6-9 than to 5 since 3,4, 6-9 are closer to origin than 5

**Cosine** emphasizes relative contributions of individual terms

# Properties of Document-Term Matrix

- Each vector  $D_i$  is a surrogate document for the original document
- Entire document-term matrix ( $N \times T$ ) is sparse, with only 0.03% of cells being non-zero in TREC collection
- Each document is a vector in terms-space
- Due to sparsity, original text documents are represented as an inverted file structure (rather than matrix directly)
  - Each term  $t_j$  points to a list of  $N$  numbers describing term occurrences for each document
- Generating document-term matrix is non-trivial
  - Are plural and singular terms counted?
  - Are very common words used as terms?

# Vector Space Representation of Queries

- Queries
  - Expressed using same term-based representation as documents
  - Query is a document with very few terms
- Vector space representation of queries
  - *database* = (1,0,0,0,0,0)
  - *SQL* = (0,1,0,0,0,0)
  - *regression* = (0,0,0,1,0,0)

Terms	
t1	database
t2	SQL
t3	index
t4	regression
t5	likelihood
t6	linear



# Query match against database using cosine distance

	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

$database = (1,0,0,0,0,0)$  → closest match is D2

$SQL = (0,1,0,0,0,0)$  → closest match is D3

$regression = (0,0,0,1,0,0)$  → closest match is D9

Use of cosine distance results in D2, D3 and D9 ranked as closest

# Weights in Vector Space Model

- $d_{ik}$  is the weight for the  $k^{\text{th}}$  term
- Many different choices for weights in IR literature
  - Boolean approach of setting weights to 1 if term occurs and 0 if it doesn't
    - favors larger documents since larger document is more likely to include query term somewhere
  - TF-IDF weighting scheme is popular
    - TF (term frequency) is same as seen earlier
    - IDF (inverse document frequency) favors terms that occur in relatively few documents

	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

# Inverse Document Frequency of a Term

- Definition

$$\log(N / n_j)$$

$N$  = total number of documents

$n_j$  = number of documents containing term  $j$

$(n_j / N)$  is the fraction of documents containing term  $j$

- IDF favors terms that occur in relatively few documents
- Example of IDF

IDF weights of terms (using natural logs):

0.105, 0.693, 0.511, 0.693, 0.357, 0.69

Term “database” occurs in many documents and is given least weight

“regression” occurs in fewest documents and given highest weight

	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

# TF-IDF Weighting of Terms

- TF-IDF weighting scheme
- TF = term frequency, denoted  $TF(d,t)$
- IDF = inverse document frequency,  $IDF(t)$
- TF-IDF weight is the product of TF and IDF for a particular term in a particular document
  - $TF(d,t)IDF(t)$
  - Example is given next

# TF-IDF Document Matrix

## TF (d,t) Document Matrix

	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

IDF (t) weights (using natural logs):  
0.105, 0.693, 0.511, 0.693, 0.357, 0.69

## TF-IDF Document Matrix

2.53	14.56	4.6	0	0	2.07
3.37	6.93	2.55	0	1.07	0
1.26	11.09	2.55	0	0	0
0.63	4.85	1.02	0	0	0
4.53	21.48	10.21	0	1.07	0
0.21	0	0	12.47	2.5	11.09
0	0	0.51	22.18	4.28	0
0.31	0	0	15.24	1.42	1.38
0.1	0	0	23.56	9.63	17.33
0.63	0	0	11.78	1.42	15.94

# Classic Approach to Matching Queries to Documents

- Represent queries as term vectors
  - 1s for terms occurring in the query and 0s everywhere else
- Represent term vectors for documents using TF-IDF for the vector components
- Use the cosine distance measure to rank the documents in terms of distance to query
  - Has disadvantage that shorter documents have a better match with query terms

# Document Retrieval with TF and TF-IDF

## TF: Document Term Matrix

	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

## TF-IDF Document Matrix

2.53	14.56	4.6	0	0	2.07
3.37	6.93	2.55	0	1.07	0
1.26	11.09	2.55	0	0	0
0.63	4.85	1.02	0	0	0
4.53	21.48	10.21	0	1.07	0
0.21	0	0	12.47	2.5	11.09
0	0	0.51	22.18	4.28	0
0.31	0	0	15.24	1.42	1.38
0.1	0	0	23.56	9.63	17.33
0.63	0	0	11.78	1.42	15.94

Query contains both *database* and *index*  
 $Q=(1,0,1,0,0,0)$

Document	TF distance	TF-IDF distance
D1	0.7	0.32
D2	0.77	0.51
D3	0.58	0.24
D4	0.6	0.23
D5	0.79	0.43
D6	0.14	0.02
D7	0.06	0.01
D8	0.02	0.02
D9	0.09	0.01
D10	0.01	0

Max values for query (1,0,1,0,0,0)

Cosine distance is high when there is a better match

TF-IDF chooses D2 while TF chooses D5  
 Unclear why D5 is a better choice  
 Cosine distance favors shorter documents  
 (a disadvantage)

# Comments on TF-IDF Method

- TF-IDF-based IR system
  - first builds an inverted index with TF and IDF information
  - Given a query (vector) lists some number of document vectors that are most similar to the query
- TF-IDF is superior in precision-recall compared to other weighting schemes
- Default baseline method for comparing performance