



Unstructured Data Processing and Information Retrieval

Week 11

COS60009: Data Management for the Big Data Age



1

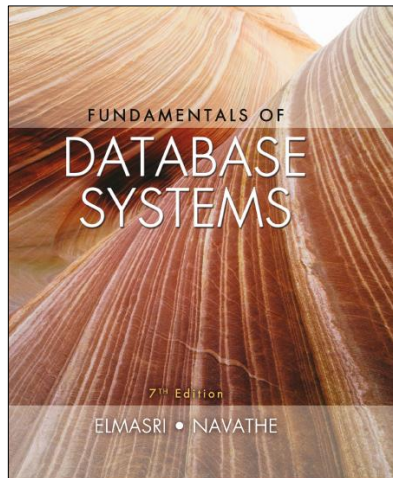
Learning Objectives

- Information Retrieval (IR) Concepts
- Retrieval Models
- Types of Queries in IR Systems
- Regular Expressions
- Text Pre-processing
- Inverted Indexing

2

Fundamentals of Database Systems

Seventh Edition



Regular Expression

Chapter 27

Introduction to
Information Retrieval and
Web Search



Copyright © 2016, 2011, 2007 Pearson Education, Inc. All Rights Reserved

3

Information Retrieval (IR) Concepts (1 of 2)

- Information retrieval
 - Process of retrieving documents from a collection in response to a query (search request)
 - Deals mainly with unstructured data
- Unstructured information
 - Does not have a well-defined formal model
 - Based on an understanding of natural language
 - Stored in a wide variety of standard formats
- Information retrieval field predates database field
 - Academic programs in Library and Information Science
- User's information need expressed as free-form search request
 - Keyword search query



Copyright © 2016, 2011, 2007 Pearson Education, Inc. All Rights Reserved

4

Information Retrieval (IR) Concepts (2 of 2)

- Characterizing an IR system
 - Types of users: Expert, Layperson
 - Types of data could be domain-specific
 - Types of information needs
 - Navigational search - looking to reach a particular website
 - Informational search - looking for a specific bit of information
 - Transactional search - getting to a website where there will be more interaction, e.g. buying something, downloading something, signing up or registering etc.
- Enterprise search systems
 - Limited to an intranet
- Desktop search engines
 - Searches an individual computer system
- IR system has no fixed data model
 - Databases have fixed schemas



Copyright © 2016, 2011, 2007 Pearson Education, Inc. All Rights Reserved

5

Comparing Databases and IR Systems

Table 27.1 A comparison of databases and IR systems

Databases	IR Systems
<ul style="list-style-type: none"> • Structured data • Schema driven • Relational (or object, hierarchical, and network) model is predominant • Structured query model • Rich metadata operations • Query returns data • Results are based on exact matching (always correct) 	<ul style="list-style-type: none"> • Unstructured data • No fixed schema; various data models (e.g., vector space model) • Free-form query models • Rich data operations • Search request returns list or pointers to documents • Results are based on approximate matching and measures of effectiveness (may be imprecise and ranked)



Copyright © 2016, 2011, 2007 Pearson Education, Inc. All Rights Reserved

6

Modes of Interactions in IR Systems

- Primary modes of interaction
 - Retrieval
 - Extract relevant information from document repository
 - Browsing
 - Exploratory activity based on user's assessment of relevance
- Web search combines both interaction modes
 - Rank of a web page measures its relevance to query that generated the result set

Generic IR Pipeline (1 of 2)

- Statistical approach
 - Documents analyzed and broken down into chunks of text
 - Each word or phrase is counted, weighted, and measured for relevance or importance
- Types of statistical approaches
 - Boolean
 - Vector space
 - Probabilistic

Generic IR Pipeline (2 of 2)

- Semantic approaches
 - Use knowledge-based retrieval techniques
 - Rely on syntactic, lexical, sentential, discourse-based, and pragmatic levels of knowledge understanding
 - Also apply some form of statistical analysis

Figure 27.1 Generic IR Framework

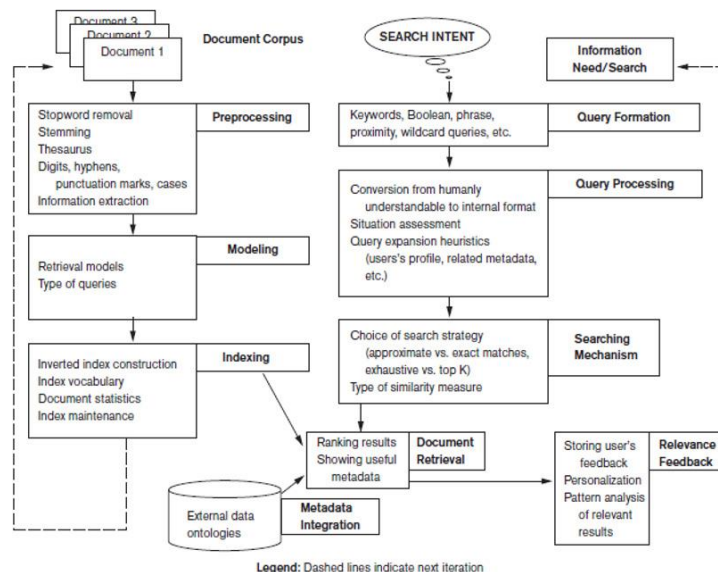
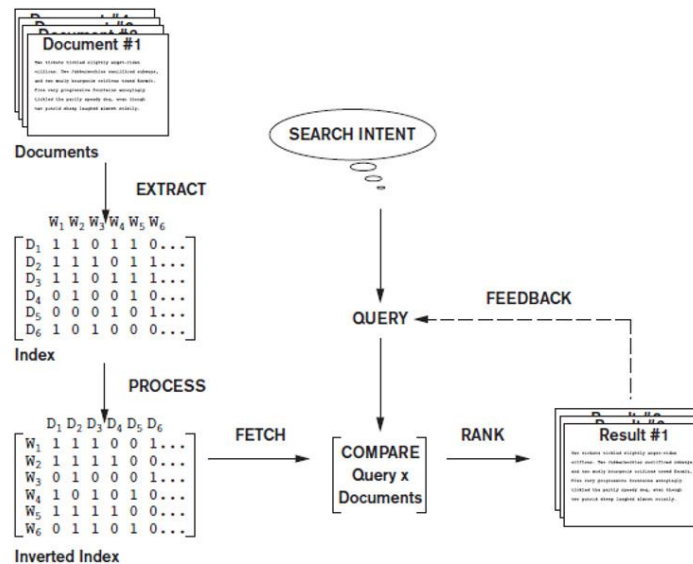


Figure 27.2 Simplified IR Process Pipeline



Retrieval Models (1 of 4)

- Boolean model
 - One of earliest and simplest IR models
 - Documents represented as a set of terms
 - Queries formulated using AND, OR, and NOT
 - Retrieved documents are an exact match
 - No notion of ranking of documents
 - Easy to associate metadata information and write queries that match contents of documents

Retrieval Models (2 of 4)

- Vector space model
 - Weighting, ranking, and determining relevance are possible
 - Uses individual terms as dimensions
 - Each document represented by an n-dimensional vector of values
 - Features
 - Subset of terms in a document set that are deemed most relevant to an IR search for the document set
 - Different similarity assessment functions can be used
- Term frequency-inverse document frequency (TF-IDF)
 - Statistical weight measure used to evaluate the importance of a document word in a collection of documents
 - A discriminating term must occur in only a few documents in the general population



Copyright © 2016, 2011, 2007 Pearson Education, Inc. All Rights Reserved

13

Retrieval Models (3 of 4)

- Probabilistic model
 - Involves ranking documents by their estimated probability of relevance with respect to the query and the document
 - IR system must decide whether a document belongs to the relevant set or nonrelevant set for a query
 - Calculate probability that document belongs to the relevant set
 - BM25: a popular ranking algorithm



Copyright © 2016, 2011, 2007 Pearson Education, Inc. All Rights Reserved

14

Retrieval Models (4 of 4)

- Semantic model
 - Morphological analysis: analyze roots and affixes to determine parts of speech of search words
 - Syntactic analysis: parse and analyze complete phrases in documents
 - Semantic analysis: resolve word ambiguities and generate relevant synonyms based on semantic relationships
 - Uses techniques from artificial intelligence and expert systems



Copyright © 2016, 2011, 2007 Pearson Education, Inc. All Rights Reserved

15

Types of Queries in IR Systems (1 of 4)

- Keyword queries
 - Simplest and most commonly used
 - Keyword terms implicitly connected by logical AND
- Boolean queries
 - Allow use of AND, OR, NOT, and other operators
 - Exact matches returned
 - No ranking possible



Copyright © 2016, 2011, 2007 Pearson Education, Inc. All Rights Reserved

16

Types of Queries in IR Systems (2 of 4)

- Phrase queries
 - Sequence of words that make up a phrase
 - Phrase enclosed in double quotes
 - Each retrieved document must contain at least one instance of the exact phrase



Copyright © 2016, 2011, 2007 Pearson Education, Inc. All Rights Reserved

17

Types of Queries in IR Systems (3 of 4)

- Proximity queries
 - How close within a record multiple search terms are to each other
 - Phrase search is most commonly used proximity query
 - Specify order of search terms
 - NEAR, ADJ (adjacent), or AFTER operators
 - Sequence of words with maximum allowed distance between them
 - Computationally expensive
 - Suitable for smaller document collections rather than the Web



Copyright © 2016, 2011, 2007 Pearson Education, Inc. All Rights Reserved

18

Types of Queries in IR Systems (4 of 4)

- Wildcard queries
 - Supports **regular expressions** and pattern-based matching
 - Example 'data*' would retrieve data, database, dataset, etc.
 - Not generally implemented by Web search engines
- Natural language queries
 - Definitions of textual terms or common facts
 - Semantic models can support

Regular Expressions

- Regular expressions are a useful way to concisely define the syntax and 'pattern' of textual data.
- Regular expressions are strings that describe the 'pattern' or 'rules' for strings
- Regular expressions are strings that follow a set of syntax rules
- Regular expressions can be used as a concise and consistent way to test for matching patterns

Regular Expressions Matching

- Perform a regular expression match `preg_match()` in PHP

- Initialise a Regular Expression pattern, and test string

```
$pattern1 = "/(chapter \d+(\.\d)*)/i";
$str1 = "For more information, see Chapter
3.4.5.1";
if (preg_match($pattern1, $str1) {
    echo "A match was found.";
} else {
    echo "A match was not found.";
}
```

Regular Expressions - Basic Examples

<code>/WebProg/</code>	matches "Isn't WebProg great?"
<code>/^WebProg/</code>	matches "WebProg rules!", not "What is WebProg?"
<code>/WebProg\$/</code>	matches "I love WebProg", not "WebProg is great!"
<code>/^WebProg\$/</code>	matches "WebProg", and nothing else
<code>/bana?na/</code>	matches "banana" and "banna", but not "banaana".
<code>/bana+na/</code>	matches "banana" and "banaana", but not "banna".
<code>/bana*na/</code>	matches "banna", "banana", and "banaaana", but not "bnana"
<code>/^[a-zA-z]+\$/</code>	matches any string of one or more letters and nothing else.

Regular Expressions – More Examples

<code>/[A-Za-z0-9-]+/</code>	letters, numbers and hyphens
<code>/\d{1,2}\d{1,2}\d{4}/</code>	date such as 19/9/2006
<code>/S+\. (jpg gif png)</code>	jpg, gif or png filename
<code>/^[1-9]{1}\$ ^[1-2]{1}[0-9]{1}\$ ^30\$/</code>	any number from 1 to 30
<code>/#?([A-Fa-f0-9]){3}((A-Fa-f0-9){3})?/</code>	valid hex colour code
<code>/(?=.*\d)(?=.*[a-z])(?=.*[A-Z]).{8,15}/</code>	password string (at least 1 uppercase letter, 1 lowercase letter and 1 digit)
<code>/^.+@.+.{2,3}\$/</code>	possible email address (in US)
<code><(/?[^\>]+)>/</code>	HTML tags

Regular Expressions - Basic Syntax

Quantifiers

<code>*</code>	0 or more
<code>+</code>	1 or more
<code>?</code>	0 or 1
<code>{4}</code>	exactly 4
<code>{4,}</code>	4 or more
<code>{4,6}</code>	4,5 or 6

Groups & Ranges

<code>.</code>	Any character (except \n)
<code>(a b)</code>	a or b
<code>(abc)</code>	group of abc (a followed by b followed by c)
<code>[abc]</code>	set ("range") a, b or c
<code>[^abc]</code>	not a, b or c
<code>[a-g]</code>	set range a to g
<code>[3-6]</code>	set range of digits 3,4,5 and 6

25

Regular Expressions - Basic Syntax

Pattern Basics

<code>^</code>	Start of string
<code>\$</code>	End of string
<code>.</code>	Match any single character
<code>(a b)</code>	a or b
<code>(...)</code>	Group section
<code>[abc]</code>	match any character in the set
<code>[^abc]</code>	not match in the set
<code>[a-z]</code>	match the range
<code>\d</code>	match a single digit from 0 to 9 <i>shortcut for [0-9]</i>
<code>\D</code>	not a digit

Pattern Quantifiers

<code>a?</code>	0 or 1 of a
<code>a*</code>	0 or more of a
<code>a+</code>	one or more instance of a
<code>a{3}</code>	exactly 3 a's = aaa
<code>a{3,}</code>	3 or more a's
<code>a{3,6}</code>	between 3 to 6 a's
<code>!(pattern)</code>	"not" pattern

[\ ^ \$. | ? * + () are the 11 **meta-characters**, or special characters, used in the syntax. If you want to include these, you need to escape them with \
eg. \ (and \.



Copyright © 2016, 2011, 2007 Pearson Education, Inc. All Rights Reserved

25

Regular Expressions - Basic Syntax

Pattern Modifiers

<code>/g</code>	global matching
<code>/i</code>	case insensitive
<code>/s</code>	single line mode
<code>/m</code>	multiple-line mode
<code>/x</code>	allow comments and white space in pattern
<code>/e</code>	evaluate replacement
<code>/U</code>	ungreedy replacement

Assertions

<code>(?=</code>	look ahead positive
<code>(?!</code>	look ahead negative
<code>(?<=</code>	look behind positive
<code>(?<!</code>	look behind negative

There are many useful online syntax references about Regular Expressions, such as:
<http://www.php.net/manual/en/reference.pcre.pattern.syntax.php>



Copyright © 2016, 2011, 2007 Pearson Education, Inc. All Rights Reserved

26

27

Regular Expressions - Basic Syntax

Anchors

<code>^</code>	start of string
<code>\A</code>	start of string
<code>\$</code>	end of string
<code>\Z</code>	end of string
<code>\b</code>	word boundary
<code>\B</code>	not word boundary
<code><</code>	start of word
<code>></code>	end of word

Special Characters

<code>\n</code>	newline
<code>\r</code>	carriage return
<code>\t</code>	tab
<code>\v</code>	vertical tab
<code>\f</code>	form feed
<code>\xhh</code>	hex character hh
<code>\s</code>	whitespace character
<code>\S</code>	not a whitespace character

Swinburne Library eBooks, can provide examples of Regular Expressions in many different scripting languages, e.g. search for 'regular expressions' in <http://proquest.safaribooksonline.com/search>



Copyright © 2016, 2011, 2007 Pearson Education, Inc. All Rights Reserved

27

Text Preprocessing (1 of 3)

- Stopword removal must be performed before indexing
- Stopwords
 - Words that are expected to occur in 80% or more of the documents of a collection
 - Examples: the, of, to, a, and, said, for, that
 - Do not contribute much to relevance
- Queries preprocessed for stopword removal before retrieval process
 - Many search engines do not remove stopwords



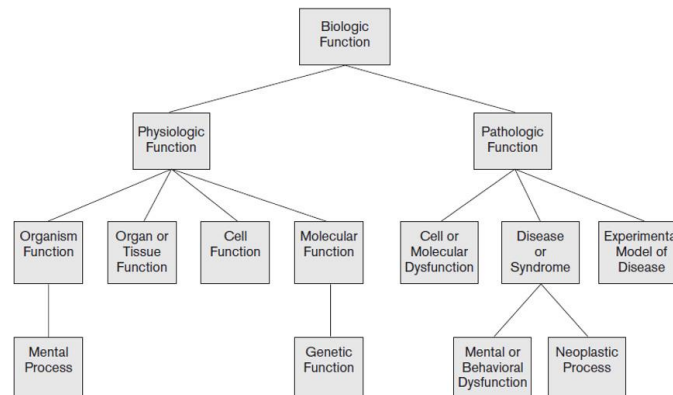
Copyright © 2016, 2011, 2007 Pearson Education, Inc. All Rights Reserved

28

Text Preprocessing (2 of 3)

- Stemming
 - Trims suffix and prefix
 - Reduces the different forms of the word to a common stem
 - Martin Porter's stemming algorithm
- Utilizing a thesaurus
 - Important concepts and main words that describe each concept for a particular knowledge domain
 - Collection of synonyms
 - UMLS

Figure 27.3 A Portion of the UMLS Semantic Network: “Biologic Function” Hierarchy



Source: UMLS Reference Manual, National Library of Medicine

Text Preprocessing (3 of 3)

- Other preprocessing steps
 - Digits
 - May or may not be removed during preprocessing
 - Hyphens and punctuation marks
 - Handled in different ways
 - Cases
 - Most search engines use case-insensitive search
- Information extraction tasks
 - Identifying noun phrases, facts, events, people, places, and relationships

Inverted Indexing (1 of 3)

- Inverted index structure
 - Vocabulary information
 - Set of distinct query terms in the document set
 - Document information
 - Data structure that attaches distinct terms with a list of all documents that contain the term

Inverted Indexing (2 of 3)

- Construction of an inverted index
 - Break documents into vocabulary terms
 - Tokenizing, cleansing, removing stopwords, stemming, and/or using a thesaurus
 - Collect document statistics
 - Store statistics in document lookup table
 - Invert the document-term stream into a term-document stream
 - Add additional information such as term frequencies, term positions, and term weights



Copyright © 2016, 2011, 2007 Pearson Education, Inc. All Rights Reserved

33

Inverted Indexing (3 of 3)

- Searching for relevant documents from an inverted index
 - Vocabulary search
 - Document information retrieval
 - Manipulation of retrieved information



Copyright © 2016, 2011, 2007 Pearson Education, Inc. All Rights Reserved

34

Figure 27.4 Example of an Inverted Index

Document 1

This example shows an example of an inverted index.

Document 2

Inverted index is a data structure for associating terms to documents.

Document 3

Stock market index is used for capturing the sentiments of the financial market.

ID	Term	Document: position
1.	example	1:2, 1:5
2.	inverted	1:8, 2:1
3.	index	1:9, 2:2, 3:3
4.	market	3:2, 3:13



Copyright © 2016, 2011, 2007 Pearson Education, Inc. All Rights Reserved

35

Copyright



This work is protected by United States copyright laws and is provided solely for the use of instructors in teaching their courses and assessing student learning. Dissemination or sale of any part of this work (including on the World Wide Web) will destroy the integrity of the work and is not permitted. The work and materials from it should never be made available to students except by instructors using the accompanying text in their classes. All recipients of this work are expected to abide by these restrictions and to honor the intended pedagogical purposes and the needs of other instructors who rely on these materials.



Copyright © 2016, 2011, 2007 Pearson Education, Inc. All Rights Reserved

36