

Marking Scheme

[P118-134 (T1)]

OCTOBER 2018 / IN - SEM (T1)**F. Y. M. TECH. (Computer Engineering) (SEMESTER - I)****COURSE NAME: Program Elective II****(Information Retrieval and Web Mining)****COURSE CODE: CSPA11184A****(PATTERN 2018)**

Q.1	a.	Lossy and lossless index compression techniques <ol style="list-style-type: none"> 1. Lossless compression: All information is preserved. <ol style="list-style-type: none"> a. What we mostly do in IR. 2. Lossy compression: Discard some information 3. Several of the preprocessing steps can be viewed as lossy compression: case folding, stop words, stemming, number elimination. 4. Dictionary Compression <ul style="list-style-type: none"> • Dictionary-as-a-String • Blocking • Front coding 5. Posting List Compression <ul style="list-style-type: none"> • Variable length encoding • Variable Byte (VB) codes 	[3+3]
	b.	Signature files in IR Characteristics : <ul style="list-style-type: none"> • Word-oriented index structures based on hashing • Low overhead (10%~20% over the text size) at the cost of forcing a sequential search over the index • Suitable for not very large texts • Inverted files outperform signature files for most applications Structure :	[2+2]

		<ul style="list-style-type: none"> • Use <i>superimposed coding</i> to create signature. • Each <i>text</i> is divided into logical blocks. • A <i>block</i> contains n distinct non-common words. • Each <i>word</i> yields “word signature”. • A <i>word signature</i> is a B-bit pattern, with m 1-bit. <ul style="list-style-type: none"> ○ Each word is divided into successive, overlapping triplets. e.g. <i>free</i> --> $\square fr, fre, ree, ee \square$ ○ Each such triplet is hashed to a bit position. • The word signatures are OR’ed to form <i>block signature</i>. • Block signatures are concatenated to form the <i>document signature</i>. <p>Example</p> <p>1 Example ($n=2, B=12, m=4$)</p> <table> <tr> <th><u>word</u></th> <th><u>signature</u></th> </tr> <tr> <td><i>free</i></td> <td>001 000 110 010</td> </tr> <tr> <td><i>text</i></td> <td>000 010 101 001</td> </tr> <tr> <td><u><i>block signature</i></u></td> <td><u>001 010 111 011</u></td> </tr> </table> <p>Search</p> <ul style="list-style-type: none"> • Use hash function to determine the m 1-bit positions. • Examine each block signature for 1’s bit positions that the signature of the search word has a 1 	<u>word</u>	<u>signature</u>	<i>free</i>	001 000 110 010	<i>text</i>	000 010 101 001	<u><i>block signature</i></u>	<u>001 010 111 011</u>	
<u>word</u>	<u>signature</u>										
<i>free</i>	001 000 110 010										
<i>text</i>	000 010 101 001										
<u><i>block signature</i></u>	<u>001 010 111 011</u>										
Q.2	a.	<p>Index construction</p> <p>The inverted index of a document collection is basically a data structure that attaches each distinctive term with a list of all documents that contains the term.</p> <p>Different techniques</p> <ol style="list-style-type: none"> 1. BSBI (Blocked Sort-Based Indexing) algorithm 2. SPIMI (Single-pass in-memory indexing) algorithm 3. Distributed indexing 4. Dynamic indexing 	[2+4]								
	b.	<p>Four challenges of unstructured data/text</p> <ul style="list-style-type: none"> • No stable document collection (spider, crawler) • Invalid document, duplication, etc. 	[1+1+1+1]								

		<ul style="list-style-type: none"> • Huge number of documents (partial collection) • Multimedia documents • Great variation of document quality • Multilingual problem • Vocabularies mismatching <ul style="list-style-type: none"> ▪ Synonymy: e.g. car v.s. automobile ▪ Polysemy: table • Queries are ambiguous, they are partial specification of user's need • Content representation may be inadequate and incomplete 	
Q.3	a.	Three models of Language Models for IR A language model is a probabilistic mechanism for generating text • Language models estimate the probability distribution of various natural language phenomena – sentences, utterances, queries Language Modeling Techniques <ul style="list-style-type: none"> • N-grams • Class-based N-grams • Probabilistic CFGs • Decision Tree 	[2+2+2]
	b.	Two types of Query expansion 1. Global Analysis: (static; of all documents in collection) <ol style="list-style-type: none"> Controlled vocabulary <ol style="list-style-type: none"> Maintained by editors (e.g., medline) Manual thesaurus <ol style="list-style-type: none"> E.g. MedLine: physician, syn: doc, doctor, MD, medico Automatically derived thesaurus <ol style="list-style-type: none"> (co-occurrence statistics) Refinements based on query log mining <ol style="list-style-type: none"> Common on the web 2. Local Analysis: (dynamic) <ol style="list-style-type: none"> Analysis of documents in result set 	[2+2]

Q.4	<p>a. Latent Semantic Indexing</p> <p>Suppose that we use the term frequency as term weights and query weights. The following document indexing rules are also used:</p> <ul style="list-style-type: none"> • stop words were not ignored • text was tokenized and lowercased • no stemming was used • terms were sorted alphabetically. <p>Solved Example</p>	[2+4]
	<p>b. Ad hoc Retrieval</p> <ul style="list-style-type: none"> • Text-based retrieval • Given a query and a corpus, find the relevant items • query: textual description of information need • corpus: a collection of textual documents • relevance: satisfaction of the user's information need • "Ad-hoc" because the number of possible queries is huge <p>Pseudo relevance feedback</p> <ul style="list-style-type: none"> • The user issues a (short, simple) query. • The search engine returns a set of documents. • User marks some docs as relevant, some as nonrelevant. • Search engine computes a new representation of the information need. Hope: better than the initial query. • Search engine runs new query and returns new results. • New results have (hopefully) better recall. • Provides a method for automatic local analysis • Pseudo-relevance feedback automates the "manual" part of true relevance feedback. • Pseudo-relevance algorithm: • Retrieve a ranked list of hits for the user's query • Assume that the top k documents are relevant. • Do relevance feedback (e.g., Rocchio) • Works very well on average • But can go horribly wrong for some queries. • Several iterations can cause <i>query drift</i>. 	[2+2]