

1 Text Search

- (a) State if the following layers belong to Search Engine Architecture, DBMS Architecture, or both.
- i Index
 - ii SQL
 - iii Buffer
 - iv Ranking Algorithm
 - v Relational Operators
- (b) What is an inverted file? Why are inverted files useful?
- (c) What is the query plan for boolean search (how do we scan and how do we join)?
- (d) Suppose we have 20 documents, and there are 300 distinct terms in all the documents. We want to find the document that is the most similar to the document “Database is fun. It is my favorite subject.” Note, the multiple choice questions in this part could have more than one answer.
- (i) If we use the vector space model, what would be the dimension of the vectors?
 - (ii) How could the dimension change if we duplicate an existing document?
 - (A) Increase
 - (B) Stay the same
 - (C) Decrease

- (iii) How could the dimension change if we add more terms to the longest document?
- (A) Increase
 - (B) Stay the same
 - (C) Decrease
- (iv) How could the dimension change if we remove a document?
- (A) Increase
 - (B) Stay the same
 - (C) Decrease
- (v) How could the dimension change if we add a term that is not in any of the documents?
- (A) Increase
 - (B) Stay the same
 - (C) Decrease
- (e) Given 100 documents, consider a term t_1 that appears 2 times in document d_1 . What is the DocTermRank of t_1 in d_1 if (assume base 10 logarithm):
- (i) T_1 appears in only d_1 ?
 - (ii) T_1 appears in 10 documents?
 - (iii) T_1 appears in all documents?
 - (iv) What property of TF-IDF does this pattern show?
- (f) Given 100 documents, consider a term t_2 that appears in 10 documents. What is the DocTermRank of t_2 in d_2 if (assume base 10 logarithm):
- (i) T_2 appears 2 times in d_2 ?
 - (ii) T_2 appears 4 times in d_2 ?
 - (iii) T_2 appears 6 times in d_2 ?
 - (iv) What property of TF-IDF does this pattern show?