

Aufgabe 1

Gegeben sei eine Datenbank, die die Dokumente d1 bis d6 enthält und diese wie folgt durch die

Indexterme t1 bis t8 repräsentiert:

d1 : {t1, t4, t6, t7}

d2 : {t2, t4, t8}

d3 : {t1, t3, t4}

d4 : {t2, t6, t7}

d5 : {t1, t4}

d6 : {t1, t3, t6}.

a) Bestimmen Sie den zu dieser Datenbank gehörigen invertierten Index („Inverted File Index“).

b) Welche Treffer liefert das Boolesche Modell für die folgenden Anfragen?

„t1 AND t4“

„(t3 OR t4) AND NOT t6“

„(t1 AND t6) OR NOT t6“

c) Formen Sie die Anfragen aus dem vorherigen Aufgabenteil jeweils zu äquivalenten Anfragen in konjunktiver und disjunktiver Normalform um.

d) Geben Sie zu jeder der folgenden Treffermengen eine (möglichst kurze) Anfrage an, die die jeweilige Treffermenge liefert.

{d2}, {d6}, {d3, d5}

Aufgabe 2

Es sollen w_1, w_2, \dots unterschiedliche Wörter sein. Die Dokumente d_1, \dots, d_5 sollen folgende Wortfolgen darstellen:

d1 : $w_5, w_1, w_9, w_3, w_8, w_2$

d2 : w_9, w_8, w_3

d3 : w_2, w_3, w_8, w_7

d4 : w_9, w_1, w_8, w_2, w_3

(a) Welche Antwortmengen ergeben sich beim Booleschen Retrieval für folgende Anfragen:

q1 : $w_2 \text{ AND } (w_8 \rightarrow w_1)$

q2 : $w_7 \text{ OR } (w_1 \rightarrow \neg w_3)$

q3 : $(w_8 \rightarrow w_2) \rightarrow (w_1 \rightarrow \neg w_3)$

(b) Glauben Sie dass die Implikation in realen Systemen üblicherweise implementiert wird?

Aufgabe 3

Schreiben Sie einen Algorithmus für das Zusammenführen zweier Postinglisten zur Auswertung einer OR-Query - analog dem AND-Algorithmus aus der Sitzung.

Aufgabe 4

Implementieren Sie folgenden ersten Teil eines sehr einfachen IR Systems.

- (a) Ein Teilprogramm soll den Shakespearertextkorpus von der Webseite in einzelne Stücke zerlegen, jedem Dokument eine ID zuordnen, den Text ausgeben und den Dokumententitel ausgeben können
- (b) Eine Methode für eine lineare Suchmöglichkeit über die Texte mit Rückgabe der Dokumenten-id
- (c) Implementieren Sie eine Klasse, die eine binäre Term-Dokumentenmatrix für einen Textkorpus wie er in a eingelesen wurde aufbaut.