

Übungen zu Multimedia-Datenbanken

Aufgabenblatt 2

Übung: Stephan Wirth und Stefan Wirtz

Vorlesung: Dr.-Ing. Marcin Grzegorzek

Fachbereich Informatik, Universität Koblenz–Landau

Ausgabe: 26.05.2009

Abgabe: 07.06.2009 per Email an stwirth@uni-koblenz.de als PDF-Anhang

Format: mmdb-blatt2-nachname1-nachname2.pdf

1 Information Retrieval - Boole'sches Modell (4 Punkte)

Gegeben seien Dokumente d_1, \dots, d_6 und die Indexmenge¹ $I = \{\text{Biaggi, Rossi, Neukirchner, MotoGP, Superbike, 2006, 2007, 2008}\}$. Nehmt folgende Zuweisung von Indextermen zu Dokumenten an:

d_1 {Biaggi,Rossi,MotoGP,2006}
 d_2 {Biaggi,2007}
 d_3 {Neukirchner,Biaggi,2008,Superbike}
 d_4 {Rossi,MotoGP,2006}
 d_5 {MotoGP,2006}
 d_6 {Neukirchner,Superbike,2008}

1. Schreibt eine Query in konjunktiver Normalform die alle Dokumente über Rossi oder Biaggi aus der MotoGP in den Jahren 2006 und 2007 zurück liefert.
2. Schreibt die Query in disjunktive Normalform um.
3. Gebt die Liste der zurückgelieferten Dokumente an.
4. Was ist der Nachteil des boole'schen Modells? Warum ist es eher Daten- als Information-Retrieval?

2 Information Retrieval - Fuzzy Modell (6 Punkte)

Gegeben seien nochmal die Dokumente d_1, \dots, d_6 von Aufgabe 1 mit der reduzierten Indexmenge $I = \{\text{Biaggi, Rossi, Neukirchner, MotoGP, Superbike}\}$. Nehmt folgende Zu-

¹Zur Info: Biaggi, Rossi und Neukirchner sind Motorradrennfahrer, MotoGP und Superbike zwei Rennserien und 2006, 2007 und 2008 sind Jahre.

weisung von Indextermen zu Dokumenten an:

- d_1 {Biaggi,Rossi,MotoGP}
- d_2 {Biaggi}
- d_3 {Neukirchner,Biaggi,Superbike}
- d_4 {Rossi,MotoGP}
- d_5 {MotoGP}
- d_6 {Neukirchner,Superbike}

1. Berechnet die Korrelationsmatrix für die Indexterme anhand des Ansatzes von Ogawa, Morita und Kobayashi. Gebt sinnvolle Zwischenschritte bei der Rechnung an.
2. Berechnet die Zugehörigkeit der Indexterme zu den Dokumenten, und stellt die Ergebnisse tabellarisch dar. Gebt sinnvolle Zwischenschritte bei der Rechnung an.
3. Berechnet die folgenden Queries anhand des Fuzzy Modells und gebt eine anhand der Relevanz sortierte Ergebnisliste an.
 - a) *Rossi and Biaggi*
 - b) *not Rossi*

3 Information Retrieval - Vektorraummodell (6 Punkte)

Gegeben seien wieder die Dokumente aus Aufgabe 2. Betrachtet die Zugehörigkeitswerte aus Aufgabe 2.2 als Termgewichte. Gegeben sei die Query *Rossi and Biaggi*. Desweiteren sei eine Query durch das Dokument d_6 spezifiziert (Ähnlichkeitssuche).

1. Berechnet die Ähnlichkeit der Dokumente zu den beiden Queries mit Hilfe der euklidischen Distanz. Gebt eine nach Relevanz sortierte Ergebnisliste an. Gebt sinnvolle Zwischenschritte bei der Rechnung an.
2. Berechnet die Ähnlichkeit der Dokumente zu den beiden Queries mit Hilfe des Kosinusmaßes. Gebt eine nach Relevanz sortierte Ergebnisliste an. Gebt sinnvolle Zwischenschritte bei der Rechnung an.
3. Was sind die Nachteile bzw. Probleme des Vektorraummodells.

4 Relevance Feedback (6 Punkte)

1. Was bedeutet Relevance Feedback?
2. Basierend auf den Dokumenten und den Indextermen aus Aufgabe 2: Ein User hat Dokument d_6 als relevant und d_4 als irrelevant eingestuft. Berechnet den neuen Anfragevektor für die Anfrage $q=Biaggi and Rossi$ mit $\alpha = 1$ und $\beta = 0.5$ mit Hilfe des Verfahrens von Rocchio. Gebt sinnvolle Zwischenschritte bei der Rechnung an.
3. Berechnet die modifizierte Anfrage mit Hilfe des euklidischen Distanzmaßes. Hat sich eine Änderung gegenüber Aufgabe 3.1 ergeben, und wenn ja, wie kann man sie deuten? Gebt sinnvolle Zwischenschritte bei der Rechnung an.

5 Bewertung von Retrieval Modellen (8 Punkte)

1. Erläutert Precision, Recall und Fall-Out. Gebt auch die jeweilige Berechnungsvorschrift an.
2. Gegeben seien Dokumente d_1, \dots, d_{20} . Bezüglich einer Anfrage q seien die Dokumente $\{d_2, d_5, d_9, d_{11}, d_{14}\}$ relevant. Zwei Systeme geben die Ergebnisliste $e_1 := \{d_2, d_4, d_5, d_9\}$ und $e_2 := \{d_2, d_3, d_5, d_6, d_8, d_9, d_{11}, d_{12}\}$. Berechnet Precision, Recall, Fall-Out.
3. Wie unterscheiden die beiden Systeme sich in ihrem Verhalten?