

Multimedia Retrieval

7 Efficient Algorithms and Data Structures

Prof. Dr. Marcin Grzegorzek

Research Group for Pattern Recognition
www.pr.informatik.uni-siegen.de

Institute for Vision and Graphics
University of Siegen, Germany

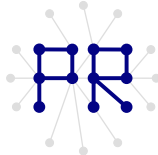


Table of Contents

Hochdimensionale
Indexstruk-
turen

Aggregation
von
Ähnlichkeiten

1 Introduction

- 1.1 Fundamental Concept
- 1.2 Search in a MMDBS
- 1.3 Applications of MMDBMS

2 Fundamentals of Information Retrieval

- 2.1 Introduction
- 2.2 Information Retrieval Models
- 2.3 Relevance Feedback
- 2.4 Evaluation of Retrieval Systems
- 2.5 User Profiles

Table of Contents

3 Fundamentals of Multimedia Retrieval

- 3.1 Characteristics of MM Management and Retrieval
- 3.2 Processing Pipeline of a Multimedia Retrieval Systems
- 3.3 Data of a Multimedia Retrieval System
- 3.4 Features
- 3.5 Applicability of Different Retrieval Models
- 3.6 Multimedia Similarity Model

4 Transforms for Feature Extraction

- 4.1 Fourier Transform
- 4.2 Wavelet Transform
- 4.3 Principal Component Analysis
- 4.4 Singular Value Decomposition

Table of Contents

5 Distance Functions

- 5.1 Properties and Classification
- 5.2 Distance Functions for Points
- 5.3 Distance Functions for Binary Data
- 5.4 Distance Functions for Sequences
- 5.5 Distance Functions for Sets

6 Similarity Measures

- 6.1 Introduction
- 6.2 Distance versus Similarity
- 6.3 Range of Similarity Measures
- 6.4 Concrete Similarity Measures
- 6.5 Aggregation of Similarity Values
- 6.6 Conversion of Distances into Similarity Values
- 6.7 Partial Similarity

Hochdimensionale
Indexstruk-
turen

Aggregation
von
Ähnlichkeiten

Table of Contents

Hochdimensionale
Indexstruk-
turen

Aggregation
von
Ähnlichkeiten

► 7 Efficient Algorithms and Data Structures

7.1 High-Dimensional Index Structures

7.2 Algorithms for Aggregation of Similarity Values

8 Query Processing

8.1 Introduction

8.2 Concepts of Query Processing

8.3 Database Model

8.4 Languages

9 Summary and Conclusions

Overview

Hochdimensionale
Indexstruk-
turen

Aggregation
von
Ähnlichkeiten

- 1 Hochdimensionale Indexstrukturen
- 2 Algorithmen zur Aggregation von Ähnlichkeitswerten

Einführung

Hochdimensionale
Indexstruk-
turen

Aggregation
von
Ähnlichkeiten

- Algorithmen und Datenstrukturen für **effiziente** Ergebnisberechnung bzgl. einer Anfrage
- Erweiterung von Verfahren klassischer DBS um Behandlung von Ähnlichkeitswerten bzw. Unähnlichkeitswerten
 - Übergang von Mengensemantik zur Listensemantik
 - hochdimensionale Indexstrukturen: zur effizienten Suche im hochdimensionalen Raum
 - Aggregation von Ähnlichkeitswerten: erforderlich etwa bei komplexen Anfragen

Overview

Hochdimensionale
Indexstruk-
turen

Aggregation
von
Ähnlichkeiten

- 1 Hochdimensionale Indexstrukturen
- 2 Algorithmen zur Aggregation von Ähnlichkeitswerten

Allgemeines

Hochdimensionale
Indexstruk-
turen

Aggregation
von
Ähnlichkeiten

- Strukturierung der Daten zur Unterstützung einer effizienten Suche
- klassische Datenstruktur in DBS: **B-Baum** und dessen Varianten
 - exakte Suche mit logarithmischem Aufwand
 - aber Einschränkung auf **eine** Dimension
 - ungeeignet zur Ähnlichkeitssuche im hochdimensionalen Raum

Anforderungen I

- Korrektheit und Vollständigkeit
- skalierbar bzgl. Dimensionsanzahl
- räumliche Ausdehnung der Objekte:
 - 0 Dimensionen: Punkt
 - 1 Dimension: Linie
 - 2 Dimensionen: Fläche
 - n Dimensionen: etwa Hyperwürfel

Allgemeines

Hochdimensionale
Indexstruk-
turen

Aggregation
von
Ähnlichkeiten

Anforderungen II

- Sucheeffizienz, also Anzahl Seitenzugriffe muss besser als bei sequentielltem Durchlauf sein
- viele Anfragearten
- effiziente Update-Operationen
- verschiedene Distanzfunktionen
- speicherplatzsparend

Anfragearten

Hochdimensionale
Indexstruk-
turen

Aggregation
von
Ähnlichkeiten

- Nächste-Nachbarsuche
- Approximative Nächste-Nachbarsuche
- Reverse-Nächste-Nachbarsuche
- Bereichssuche
- Punktsuche
- Partial-Match-Suche
- Ähnlichkeitsverbund

Anfragearten

Hochdimensionale
Indexstruk-
turen

Aggregation
von
Ähnlichkeiten

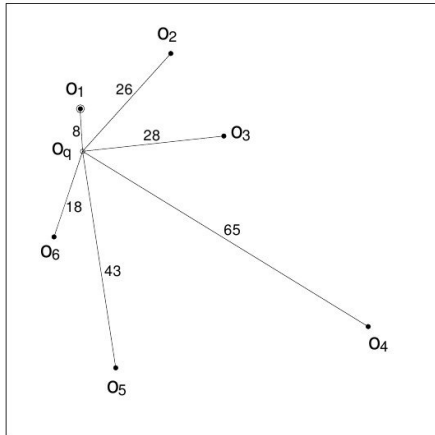
Nächste-Nachbarsuche

- Feature-Daten eines Anfragemedienobjekts: o_q
- Menge von Feature-Daten: FO
- binäre Distanzfunktion $d()$
- Finden des ähnlichsten Medienobjekts (das nächste Feature-Objekt)
- mehrere nächste Nachbarn möglich:

$$nn(o_q) \subseteq FO \text{ mit } \forall o_i \in FO : \forall o \in nn(o_q) : d(o_q, o) \leq d(o_q, o_i)$$

Anfragearten

Nächste-Nachbarsuche graphisch



Hochdimensionale
Indexstruk-
turen

Aggregation
von
Ähnlichkeiten

Anfragearten

Hochdimensionale
Indexstruk-
turen

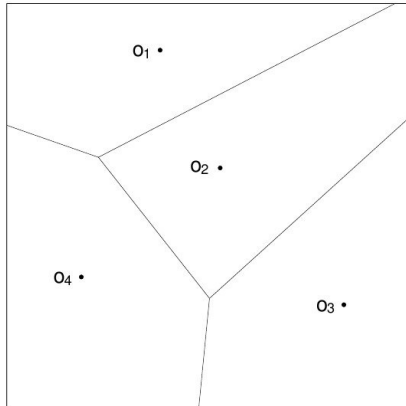
Aggregation
von
Ähnlichkeiten

Zweidimensionale Voronoi-Zellen

- NN-Suche auf punktförmigen Feature-Daten ist äquivalent zum **Enthaltsenseintest in Voronoi-Zelle**
- jedem Feature-Objekt ist eigene Voronoi-Zelle zugewiesen
- Voronoi-Zelle enthält alle Raumpunkte, die nächste Nachbarn des entspr. Feature-Objekts sind
- Idee: Vorausberechnung aller Voronoi-Zellen und anschließend Enthaltenseintest
Problem: hohe Berechnungskomplexität für Enthaltenseintest

Anfragearten

Voronoi-Zellen graphisch



Hochdimensionale
Indexstruk-
turen

Aggregation
von
Ähnlichkeiten

Anfragearten

KNN-Suche

Hochdimensionale
Indexstruk-
turen

Aggregation
von
Ähnlichkeiten

- die k nächsten Nachbarn werden gesucht

$$knn(o_q) \subseteq FO \text{ mit } |knn(o_q)| = k \text{ und } \forall o \in FO \setminus knn(o_q) : \\ \forall o_{knn} \in knn(o_q) : d(o_q, o) \geq d(o_q, o_{knn})$$

- bei gleichen Distanzen: **nichtdeterministische** Auswahl
- Ergebnisobjekte werden aufsteigend nach Distanz **sortiert**
- k ist üblicherweise so klein, das Ergebnis in Hauptspeicher passt
- ansonsten: Ergebnisobjekt sukzessive abholen
(**getNext**-Semantik/**ranking**-Anfrage)

Anfragearten

Hochdimensionale
Indexstruk-
turen

Aggregation
von
Ähnlichkeiten

Approximative Nächste-Nachbarsuche

- Effizienzgewinn bzgl. NN-Anfragen, wenn kleine Ungenauigkeiten tolerierbar
- ϵ als Maß der Ungenauigkeit

$$ann(o_q) = o \in FO \text{ wenn } d(o_q, o) \leq (1 + \epsilon) \cdot d(o_q, nn(o_q))$$

- mehrere Feature-Objekte können Bedingung erfüllen
→ nichtdeterministische Auswahl
- Vorsicht: *ann-Ergebnis muss nicht in Nähe des nn-Ergebnisses liegen*

Anfragearten

Hochdimensionale
Indexstruk-
turen

Aggregation
von
Ähnlichkeiten

PAC-NN-Suche

(probably approximative correct)

- weitere **Abschwächung** einer ANN-Suche
- Forderung: Wahrscheinlichkeit der Abweichung von *ann*-Bedingung muss vorgebbaren Mindestwert überschreiten
- ermöglicht effizientere Suche

Anfragearten

Hochdimensionale
Indexstruk-
turen

Aggregation
von
Ähnlichkeiten

Reverse-Nächste-Nachbarsuche

- Suche nach Feature-Objekten, deren nächster Nachbar der Anfragepunkt ist
(etwa Suche nach bestem Ort für neuen Einkaufsmarkt)

$$rnn(o_q) = \{o \in FO | o_q \in nn(o)\}$$

- Achtung: Ergebnis oft anderes als bei NN-Suche, da
Nächste-Nachbarrelation **nicht symmetrisch** ist



Anfragearten

Hochdimensionale
Indexstruk-
turen

Aggregation
von
Ähnlichkeiten

Bereichssuche

- Anfrage definiert einen **Bereich** (**Region**) im hochdimensionalen Raum
- Ergebnis sind alle Feature-Objekte, die Anfragebereich schneiden

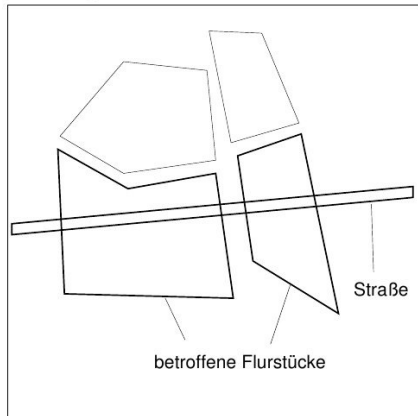
$$\text{range}(o_q) = \{o \in FO \mid o \cap o_q \neq \emptyset\}$$

- Varianten
 - begrenzte versus unbegrenzte Bereiche
 - Spezialfall **Hyperkugel** und **Hyperrechteck**

Anfragearten

Bereichssuche graphisch

Beispiel: Straßenplanung im Katasteramt:



Hochdimensionale
Indexstruk-
turen

Aggregation
von
Ähnlichkeiten

Anfragearten

Hochdimensionale
Indexstruk-
turen

Aggregation
von
Ähnlichkeiten

Punktsuche

- Suche anhand eines gegebenen Feature-Objekts o_q
- Test auf **Enthaltensein** (exakte Überdeckung)

$$\text{punkt}(o_q) = \begin{cases} \text{wahr} & : \exists o \in FO : o_q = o \\ \text{falsch} & : \text{sonst} \end{cases}$$

- Punktsuche in MMDB ist **relativ selten**

Anfragearten

Hochdimensionale
Indexstruk-
turen

Aggregation
von
Ähnlichkeiten

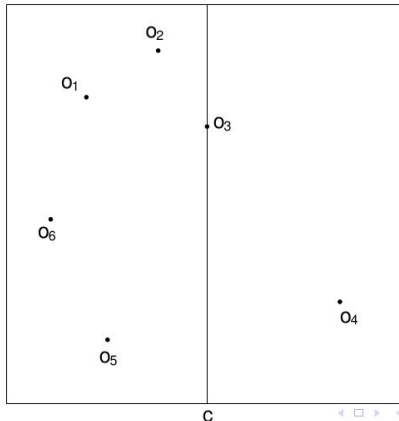
Partial-Match-Suche

- Punktsuche kann als Complete-Match-Suche aufgefasst werden
- bei Partial-Match-Suche Übereinstimmung nur in **einigen Dimensionen**
(restliche Dimensionen werden ignoriert)
- ist Spezialfall der Bereichsanfrage mit teilweise unbegrenztem Bereich

Anfragearten

Partial-Match-Suche graphisch

Suchbereich ist senkrechte Linie (Übereinstimmung in nur einer Dimension)



Hochdimensionale
Indexstruk-
turen

Aggregation
von
Ähnlichkeiten

Anfragearten

Hochdimensionale
Indexstruk-
turen

Aggregation
von
Ähnlichkeiten

Ähnlichkeitsverbund

- Operation auf **zwei Mengen** von Feature-Objekten
- **Verbund** findet Paare, deren Distanz kleiner als vorgegebener Schwellenwert ϵ ist

$$sj(FO_1, FO_2, \epsilon) = \{(o_1, o_2) | o_1 \in FO_1 \wedge o_2 \in FO_2 \wedge d(o_1, o_2) \leq \epsilon\}$$

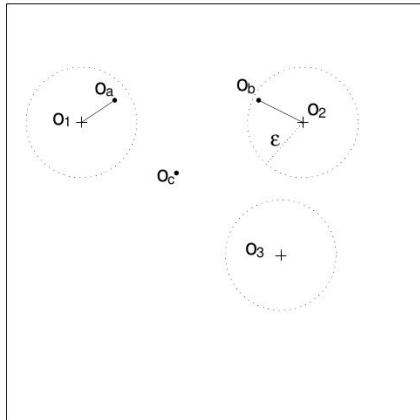
- **Selbstverbund**: dieselbe Menge zweimal

Anfragearten

Ähnlichkeitsverbund graphisch

Hochdimensionale
Indexstruk-
turen

Aggregation
von
Ähnlichkeiten



Baumverfahren

Hochdimensionale
Indexstruk-
turen

Aggregation
von
Ähnlichkeiten

- Ausgangspunkt: **punktförmige hochdimensionale** Feature-Objekte
- B-Baum ist **eindimensional**
- Abbildung mehrdimensionaler Raum auf eine Dimension im Allgemeinen nicht distanzerhaltend möglich
(siehe etwa Simplex mit $n + 1$ Punkten im n -dimensionalen Raum)
- Fazit: **mehr**dimensionale Indexverfahren sind erforderlich

Baumverfahren

Hochdimensionale
Indexstruk-
turen

Aggregation
von
Ähnlichkeiten

Grundidee hierarchischer Indexierungsverfahren:

- Beschreiben von Punktmengen durch geometrische, umschreibende **Regionen (Cluster)**
- bei der Suche Test und evtl. Ausschluss von Regionen
- Regionen können sich gegenseitig enthalten
→ **Halbordnung** und daher Hierarchie (Hasse-Diagramm)
oder Baum

Unterscheidungskriterien von Baumverfahren

Hochdimensionale
Indexstruk-
turen

Aggregation
von
Ähnlichkeiten

Merkmal	Unterscheidung
Cluster-Bildung	global zerlegend (space partitioning) lokal gruppierend (data partitioning)
Cluster-Überlappung	überlappend disjunkt
Balance	balanciert unbalanciert
Objektspeicherung	Blätter und Knoten Blätter
Geometrie	Hyperkugel Hyperquader Hyperellipsoid

Suchalgorithmen in Bäumen

Hochdimensionale
Indexstruk-
turen

Aggregation
von
Ähnlichkeiten

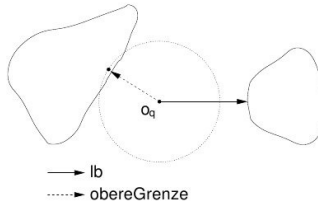
- Algorithmen zur Berechnung des **nächsten Nachbars**
- Anfragepunkt q
- Existenz zweier Distanzfunktionen
 - Distanz zwischen zwei Punkten
 - minimale Distanz zwischen q und potenziellem Clusterpunkt eines Clusters

Branch-und-Bound-Algorithmus

Hochdimensionale
Indexstruk-
turen

Aggregation
von
Ähnlichkeiten

- 1988 publiziert von Niemann und Goppert
- geht von Objektspeicherung in Blättern aus
- realisiert **Tiefensuche**
- verwendet **dynamisch** angepasste Distanz obereGrenze zu NN-Kandidat



Branch-und-Bound-Algorithmus

Hochdimensionale
Indexstruk-
turen

Aggregation
von
Ähnlichkeiten

```
[1]  real obereGrenze =  $\infty$ 
[2]  punkt naechsterNachbar = nil
[3]
[4]  procedure BranchAndBound(punkt q,knoten T)
[5]
[6]      sortiere Subknoten von T aufsteigend nach lb-Distanz zu q
[7]      for each Subknoten k von T do
[8]          if k ist Blatt then
[9]              for each Punkt p in k do
[10]                  distanz = d(p,q) //Distanz zwischen p und q
[11]                  if distanz < obereGrenze then do
[12]                      obereGrenze = distanz
[13]                      naechsterNachbar = p // NN-Kandidat
[14]                  end if
[15]              end for
[16]          else do
[17]              lb=lb(q,k) //kleinstmögliche Distanz von q zu k
[18]              if lb > obereGrenze then
[19]                  schließe k von allen weiteren Betrachtungen aus
[20]              else BranchAndBound(q,k)
[21]              end else
[22]          end for
[23] end procedure
```

RKV-Algorithmus

Hochdimensionale
Indexstruk-
turen

Aggregation
von
Ähnlichkeiten

Roussopoulos/Kelly/Vincent 1995

- spezieller Branch-and-Bound-Algorithmus zur Nächste-Nachbarsuche
- Feature-Objekte als hochdimensionale Objekte
- für Baum mit **lokal gruppierenden Hyperquadern** und Feature-Objekten in den Blättern
- Hyperquader als **MBR** (engl. minimum bounding rectangle):
 - jede Hyperfläche berührt mind. ein Feature-Objekt von außen
 - MBR wird durch zwei Punkte (s, t) (innerste und äußerste Ecke) identifiziert

RKV-Algorithmus

Hochdimensionale
Indexstruk-
turen

Aggregation
von
Ähnlichkeiten

MIN-Distanz

minimal mögliche Distanz zwischen Anfragepunkt q und MBR (s, t)

$$MINDIST(q, (s, t)) = \sum_{i=1}^n |q[i] - r[i]|^2$$

mit

$$r[i] = \begin{cases} s[i] & \text{wenn } q[i] < s[i] \\ t[i] & \text{wenn } q[i] > t[i] \\ q[i] & \text{sonst.} \end{cases}$$

RKV-Algorithmus

MINMAX-Distanz

maximal möglicher Abstand zum nächsten Nachbarn in (s, t)
(Ausnutzen der **Minimaleigenschaft** eines MBR)

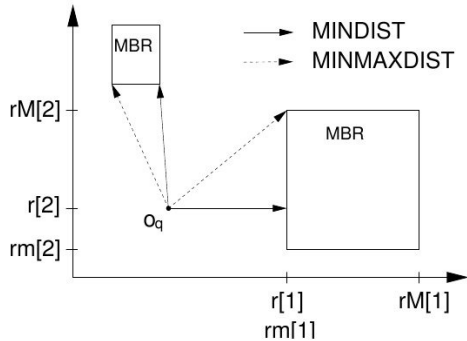
$$\text{MINMAXDIST}(q, (s, t)) = \min_{1 \leq k \leq n} \left(|q[k] - rm[k]|^2 + \sum_{\substack{i \neq k \\ 1 \leq i \leq n}} |q[i] - rM[i]|^2 \right)$$

mit

$$rm[k] = \begin{cases} s[k] & \text{wenn } q[k] \leq \frac{(s[k] + t[k])}{2} \\ t[k] & \text{sonst} \end{cases} \quad \text{und}$$
$$rM[i] = \begin{cases} s[i] & \text{wenn } q[i] \geq \frac{(s[i] + t[i])}{2} \\ t[i] & \text{sonst} \end{cases}$$

RKV-Algorithmus

MIN- und MINMAX-Distanzen graphisch



Hochdimensionale
Indexstruk-
turen

Aggregation
von
Ähnlichkeiten

RKV-Algorithmus

Hochdimensionale
Indexstruk-
turen

Aggregation
von
Ähnlichkeiten

Reduzierung des Suchaufwands

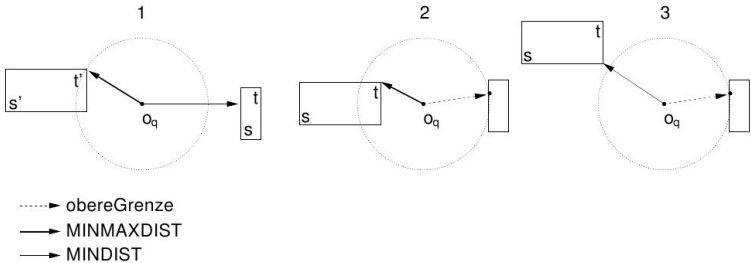
- Sortierung der Kindsnoten anhand MIN-Distanz (optimistisch) oder MINMAX-Distanz (pessimistisch)
- 3 Strategien zur Suchaufwandreduzierung (*obereGrenze* ist Distanz zum NN-Kandidaten)
 - 1 $MINDIST(q, (s, t)) > MINMAXDIST(q, (s', t')) :$
→ MBR (s, t) braucht nicht aufgesucht zu werden
 - 2 $obereGrenze > MINMAXDIST(q, (s, t)) :$
→ $obereGrenze = MINMAXDIST(q, (s, t))$
 - 3 $obereGrenze < MINDIST(q, (s, t)) :$
→ MBR (s, t) braucht nicht aufgesucht zu werden

RKV-Algorithmus

Hochdimensionale
Indexstruk-
turen

Aggregation
von
Ähnlichkeiten

Reduzierung des Suchaufwands graphisch



RKV-Algorithmus

Hochdimensionale
Indexstruk-
turen

Aggregation
von
Ähnlichkeiten

```
[1]  procedure RKV(punkt q,knoten T,real obereGrenze,
[2]                objekt naechsterNachbar)
[3]    knoten neuerKnoten
[4]    brancharray branchList
[5]    real distanz
[6]    objekt o
[7]    if T ist Blattknoten then
[8]        for each o in T do
[9]            distanz ist Distanz zwischen q und o
[10]           if distanz < obereGrenze then do
[11]               obereGrenze = distanz
[12]               naechsterNachbar = o
[13]           end if
[14]       end for
[15]    else do
[16]        branchList sei Liste von Kindknoten aus T
[17]        branchList nach MINDIST oder MINMAXDIST sortieren
[18]        branchList nach Strategie 1 kürzen
[19]        obereGrenze nach Strategie 2 reduzieren
[20]        branchList nach Strategie 3 kürzen
[21]        for each neuerKnoten in branchList do
[22]            RKV(q,neuerKnoten,obereGrenze,naechsterNachbar)
[23]            branchList nach Strategie 3 kürzen
[24]        end for
[25]    end do
[26] end procedure
```


RKV-Algorithmus

Hochdimensionale
Indexstruk-
turen

Aggregation
von
Ähnlichkeiten

knn-Anfragen mit RKV-Algorithmus

Modifikation des Algorithmus:

- **sortierte Warteschlange** zur Verwaltung der k Nächste-Nachbarnkandidaten
- *obereGrenze* ist Distanz zum **letzten** Kandidat

RKV-Algorithmus

Hochdimensionale
Indexstruk-
turen

Aggregation
von
Ähnlichkeiten

Bewertung RKV-Algorithmus

- getNext-Anfragen werden nicht unterstützt
(Problem aufgrund der Tiefensuche)
- ursprünglich für euklidische Distanz entwickelt;
funktioniert auch auf anderen Distanzfunktionen, so lange
MIN- und MINMAX-Distanzen berechnet werden können
- Einschränkung auf MBRs als Clustergeometrien
→ ansonsten statt MINMAX- die MAX-Distanz verwenden

HS-Algorithmus

Hochdimensionale
Indexstruk-
turen

Aggregation
von
Ähnlichkeiten

Henrich/Hjaltason/Samet

- Algorithmus für **getNext**-Anfragen
- statt Tiefensuche Verwendung einer **global sortierten Warteschlange**:
 - enthält Knoten, Blätter und Feature-Objekte mit minimalen Distanzen $1b$ zum Anfragepunkt
 - legt Abarbeitungsreihenfolge aufgrund aufsteigender Distanz fest
 - nur Kopfelemente werden entnommen

HS-Algorithmus

Hochdimensionale
Indexstruk-
turen

Aggregation
von
Ähnlichkeiten

- Initialisierung mit Wurzelement
- wenn entnommenes Kopfelement Knoten, dann werden dessen Kinder eingefügt
- wenn entnommenes Kopfelement Blatt, dann werden dessen Feature-Objekte eingefügt
- wenn entnommenes Kopfelement **Feature-Objekt** → Ausgabe

HS-Algorithmus

```
[1]  procedure HS(punkt q,knoten T)
[2]      pqueue queue // Priority Queue
[3]      queueEintrag element // Priority-Queue-Eintrag
[4]      objekt fo // Feature-Objekt
[5]      knoten k
[6]      enqueue(queue, T, lb(q,T))
[7]      while not isEmpty(queue) do
[8]          element = dequeue(queue)
[9]          if element ist Feature-Objekt then do
[10]             while element = first(queue) do
[11]                 deleteFirst(queue) // Duplikate entfernen
[12]             end do
[13]             ausgabe(element) // getNext-Resultat ausgeben
[14]          end do
[15]          else if element ist Blattknoten then
[16]              for each fo in element do
[17]                  enqueue(queue, fo, lb(q,fo))
[18]              end do
[19]          else // innerer Knoten
[20]              for each Kind k in element do
[21]                  enqueue(queue, k, lb(q,k))
[22]              end do
[23]          end if
[24]      end do
[25]  end procedure
```

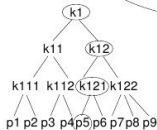
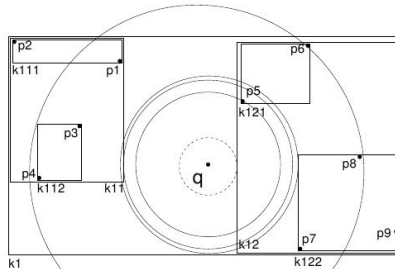
Hochdimensionale
Indexstruk-
turen

Aggregation
von
Ähnlichkeiten

HS-Algorithmus

Hochdimensionale
Indexstruk-
turen

Aggregation
von
Ähnlichkeiten



Queue: k1
k12 k11
k121 k11 k122
p5 k11 k122 p6

HS-Algorithmus

Hochdimensionale
Indexstruk-
turen

Aggregation
von
Ähnlichkeiten

Bewertung

- gut geeignet für getNext-Anfragen
- unabhängig von Cluster-Geometrie
- Problem: Warteschlange kann zu groß für Hauptspeicher werden
 - aufwändige Auslagerung auf Festplatte notwendig
- Navigation „springt“ aufgrund Warteschlangensortierung keine Tiefen/Breitensuche
 - teure Festplattenzugriffe

Overview

Hochdimensionale
Indexstruk-
turen

Aggregation
von
Ähnlichkeiten

- 1 Hochdimensionale Indexstrukturen
- 2 Algorithmen zur Aggregation von Ähnlichkeitswerten

Aggregation - Einführendes Beispiel

Hochdimensionale
Indexstruk-
turen

Aggregation
von
Ähnlichkeiten

- Gesucht sind alle Bilder, die zu einem vorgegebenen Photo bezüglich Farbverteilung und Textur ähnlich sind. Für jedes Ergebnisbild müssen also zwei Ähnlichkeitswerte anhand einer geeigneten Aggregatfunktion kombiniert werden.
- Gesucht sind alle Bilder, die zu mehreren Anfragebildern ähnlich sind. In diesem Fall müssen pro Ergebnisbild mehrere Ähnlichkeitswerte aggregiert werden.

Anforderungen

An eine Aggregation agg , die Ähnlichkeitswerte für ein Objekt aggregiert, werden bestimmte Forderungen gestellt:

Hochdimensionale
Indexstruk-
turen

Aggregation
von
Ähnlichkeiten

1. *Ähnlichkeitswerte*. Die Funktion muss mehrere Ähnlichkeitswerte aus dem Intervall $[0, 1]$ auf einen Wert aus dem Intervall $[0, 1]$ abbilden:

$$\text{agg} : [0, 1]^n \rightarrow [0, 1]$$

2. *Monotonie*. Wenn die Eingangswerte nicht sinken, dann sinkt auch der aggregierte Ähnlichkeitswert nicht:

$$x_1 \leq y_1 \wedge \dots \wedge x_n \leq y_n \Rightarrow \text{agg}(x_1, \dots, x_n) \leq \text{agg}(y_1, \dots, y_n)$$

Anforderungen

Hochdimensionale
Indexstruk-
turen

Aggregation
von
Ähnlichkeiten

3. *Strikte Monotonie.* Wenn alle Eingangswerte wachsen, dann muss auch der entsprechende, aggregierte Ähnlichkeitswert wachsen:

$$x_1 < y_1 \wedge \dots \wedge x_n < y_n \Rightarrow \text{agg}(x_1, \dots, x_n) < \text{agg}(y_1, \dots, y_n)$$

4. *Stetigkeit.* Die Aggregatfunktion soll bezüglich der Eingangswerte stetig sein, also keine abrupten Sprünge aufweisen.

Anforderungen

Hochdimensionale
Indexstruk-
turen

Aggregation
von
Ähnlichkeiten

5. *Idempotenz*. Eine Aggregation derselben Werte muss diesen Wert selbst ergeben:

$$\text{agg}(a, \dots, a) = a$$

6. *Unabhängigkeit von der Reihenfolge*. Das Resultat einer Aggregation ist unabhängig von der Reihenfolge der zu aggregierenden Ähnlichkeitswerte:

$$\text{agg}(x_1, x_2, \dots, x_n) = \text{agg}(x_{p_1}, x_{p_2}, \dots, x_{p_n})$$

wobei $[p_i]$ eine beliebige Permutation der Werte $[i]$ darstellt.

Generalisiertes Mittel

Hochdimensionale
Indexstruk-
turen

Aggregation
von
Ähnlichkeiten

- Das generalisierte Mittel ist definiert wie folgt:

$$\text{agg}_{gm}^{\alpha}(x_1, \dots, x_n) = \left(\frac{x_1^{\alpha} + \dots + x_n^{\alpha}}{n} \right)^{\frac{1}{\alpha}}$$

- Der Parameterwert α muss ungleich 0 sein.
- Folgende Spezialfälle ergeben sich:
 - $\alpha = 1$: arithmetisches Mittel
 - $\alpha = \infty$: maximaler Ähnlichkeitswert
 - $\alpha = -\infty$: minimaler Ähnlichkeitswert

Klassifikation

Hochdimensionale
Indexstruk-
turen

Aggregation
von
Ähnlichkeiten

1. *Combiner-Algorithmen*: Ausgangspunkt dieser Algorithmen sind mehrere, nach Ähnlichkeitswerten absteigend sortierte Objektlisten, die anhand einer Aggregatfunktion zu einer solchen sortierten Liste vereinigt werden sollen.
2. *Kondensator-Algorithmen*: Ausgangspunkt für diese Algorithmen ist eine Liste von Objekten, bei der mehrere Listenobjekte zu jeweils einem neuen Listenobjekt aggregiert werden.
3. *Indexaggregation*: Bei diesen Algorithmen existieren keine Eingangslisten. Statt dessen wird innerhalb einer Indexstruktur die Aggregation ausgeführt, bevor eine sortierte Liste erstellt wird.

Combiner-Algorithmen

Hochdimensionale
Indexstruk-
turen

Aggregation
von
Ähnlichkeiten

- Jeder Listeneintrag `e` enthält einen eindeutigen Identifikator `e.id` eines Medienobjekts zusammen mit einem Ähnlichkeitswerte `e.grade`.
- Wir beschränken uns der Einfachheit halber auf zwei Eingangslisten: `links` und `rechts`. Ein Objekt `o` hat einen Eintrag `o.lgrade` für den Ähnlichkeitswert der linken Liste und `o.rgrade` der rechten Liste.
- Der aggregierte Ähnlichkeitswert wird mit `o.grade` und der Identifikator mit `o.id` notiert. Einträge, denen noch kein Wert zugewiesen wurde, besitzen den Default-Wert `NULL`.

Combiner-Algorithmen

Hochdimensionale
Indexstruk-
turen

Aggregation
von
Ähnlichkeiten

- Auf die Listenelemente kann in zwei verschiedenen Modi zugegriffen werden. Zum einen ist ein sequentieller Zugriff, also jeweils ein Zugriff auf den Beginn der Liste, möglich. Zusätzlich kann auf ein Medienobjekt über den dazugehörenden Identifikator direkt zugegriffen werden.
- Im Folgenden werden die beiden Zugriffsmodi mit *sequentiell*em Zugriff und mit *randomisiert*em Zugriff bezeichnet. Für die Zugriffe werden die Funktionen `getNext()` und `random()` verwendet.

Combiner-Algorithmen

Hochdimensionale
Indexstruk-
turen

Aggregation
von
Ähnlichkeiten

- Die Ähnlichkeitswerte eines Medienobjekts aus den verschiedenen Listen sollen anhand einer monotonen Aggregatfunktion agg zu einem neuen Ähnlichkeitswert zusammengefasst werden.
- Von allen Ergebnisobjekten soll auf k Objekte mit den größten, aggregierten Ähnlichkeitswerten effizient zugegriffen werden können.
- Alternativ wird oft auch ein sequentieller Zugriff auf die Objekte in der absteigenden Reihenfolge der aggregierten Ähnlichkeitswerte verlangt.
- Die beiden Arten der Ergebniszugriffe werden mit $\text{top-}k$ - und mit ranking- Zugriff bezeichnet.

Combiner-Algorithmen

Hochdimensionale
Indexstruk-
turen

Aggregation
von
Ähnlichkeiten

- TA-Algorithmus
- NRA-Algorithmus
- Stream-Combine-Algorithmus

Combiner-Algorithmen: TA-Algorithmus

Hochdimensionale
Indexstruk-
turen

Aggregation
von
Ähnlichkeiten

```
01 procedure TA(liste links, liste rechts, funktion agg, liste top-k)
02   eintrag ol, or
03   real tau
04   repeat
05     ol = getNext(links)
06     or = getNext(rechts)
07     ol.rgrade = random(rechts, ol.id)
08     or.lgrade = random(links, or.id)
09     ol.grade = agg(ol.lgrade, ol.rgrade)
10     or.grade = agg(or.lgrade, or.rgrade)
11     aktualisiere top - k bzgl. ol und or
12     tau = agg(ol.lgrade, or.rgrade)
13   until  $|top - k| = k$  and  $\forall o \in top - k : o.grade \geq tau$ 
14   sortier top - k - Elementenach grade
15 end procedure
```

Combiner-Algorithmen: TA-Algorithmus

Hochdimensionale
Indexstruk-
turen

Aggregation
von
Ähnlichkeiten

Beispiel

links		rechts		tau	top-k	
id	lgrade	id	rgrade	tau	id	lgrade+rgrade
o3	0,9	o4	0,8	1,7	o4	1,4
o1	0,7	o2	0,7	1,4	o1	1,3
o4	0,6	o1	0,6	1,2	–	–
o2	0,2	o5	0,4	0,6	–	–
o5	0,1	o3	0,1	0,2	–	–