

Free text of gecontroleerd vocabulaire: een dilemma

Dr. Gerhard J.A. Riesthuis
Universiteit van Amsterdam
17 februari 2004



17 feb. 2004

dr. Gerhard J.A. Riesthuis

1

Een oude discussie

- | In de jaren dertig van de 20ste eeuw
 - Is onderwerpsontsluiting nodig?
 - Is een trefwoordencatalogus beter dan een systematische catalogus?



17 feb. 2004

dr. Gerhard J.A. Riesthuis

2

Onderwerpsontsluiting nodig?

Neen

- | *Het is zeer duur*

- | *Het is overbodig*
 - *Het kan ook met titelwoorden*
 - *Er bestaan vele bibliografiën*



17 feb. 2004

dr. Gerhard J.A. Riesthuis

3

Onderwerpsontsluiting nodig?

Ja

- | De omweg via bibliografiën is te omslachtig en verschijnen te laat
- | Bibliografiën bevatten niet altijd ook boeken
- | Systematische ontsluiting nodig voor overzicht van collecties
- | Bibliotheek kan onderwerpsontsluiting afstemmen op eigen gebruikers



17 feb. 2004

dr. Gerhard J.A. Riesthuis

4

RESULTAAT VAN DISCUSSIE

- | Marburg stopte met onderwerpsontsluiting
- | De andere UB's gingen door met een onderwerpscatalogus
- | Wel verschuiving van systematische naar trefwoordcatalogi



17 feb. 2004

dr. Gerhard J.A. Riesthuis

5

Waarom die discussie toen?

Documentatiebeweging is volwassen

- | Veel bibliografieën beschikbaar
- | Verschuiving van boeken naar tijdschriftartikelen
- | Voor literatuuronderzoek studie van bibliografieën nodig (voor artikelen)



17 feb. 2004

dr. Gerhard J.A. Riesthuis

6

En toen?

- | 50 jaar later stopte ook de Universiteitsbibliotheek van Gent met onderwerpsontsluiting
- | Argumenten
 - Bibliografieën
 - Titelwoorden
- | Maar: besluit wordt nu betreurd



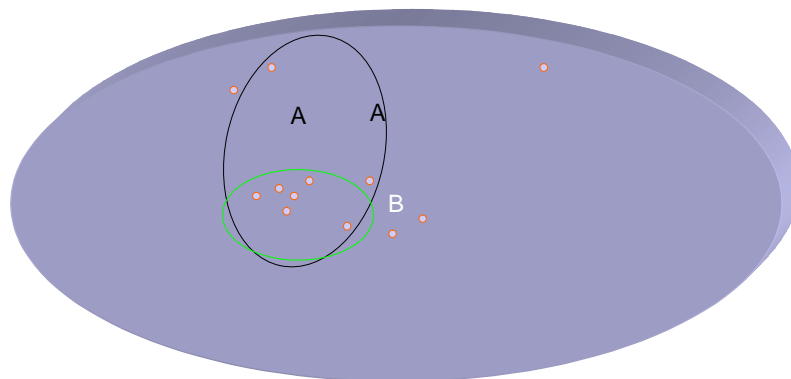
17 feb. 2004

dr. Gerhard J.A. Riesthuis

7

Hoe werkt zoeken?

Collectie met gezochte documenten.



17 feb. 2004

dr. Gerhard J.A. Riesthuis

8

In woorden:

- | Zoeken betekent het vormen van een deelverzameling waarin het aandeel van relevante documenten groter is dan in de verzameling als geheel
- | Zonder te veel relevante documenten te verliezen
- | Ideaal: deelverzameling met alle relevante en geen niet-relevante documenten



Vele wegen naar het doel

- | Onderwerpsontsluiting met gecontroleerde vocabulaires
- | Titelwoorden
- | Vrije tekst [tekst digitaal beschikbaar!]
- | Andere
 - Namen van uitgevers
 - Titels van reeksen
 - Namen van auteurs, redacteurs, enz.



Kosten

- | Relatief duur: Onderwerpsontsluiting
- | Relatief goedkoop: De overige methoden
- | Vragen:
 - Zijn de resultaten van onderwerpsontsluiting zoveel beter dat de extra kosten verantwoord zijn?
 - Kunnen de kosten verlaagd worden zonder (veel) kwaliteitsverlies?



Gecontroleerde vocabulaires

- | Classificaties
 - Enumeratieve classificaties
 - Facetclassificaties (analytisch-synthetisch)
- | Woordsystemen
 - Trefwoordsystemen
 - Thesauri



Classificaties à β Woordsystemen

I Classificaties

- geven overzicht
- relatief moeilijk bij het zoeken
- relatief betere resultaten bij ontsluiten

I Woordsystemen

- relatief makkelijk bij het zoeken,
- goede resultaten bij zoeken op “namen”, maar slecht bij vage onderwerpen
- relatief slechtere resultaten bij ontsluiten

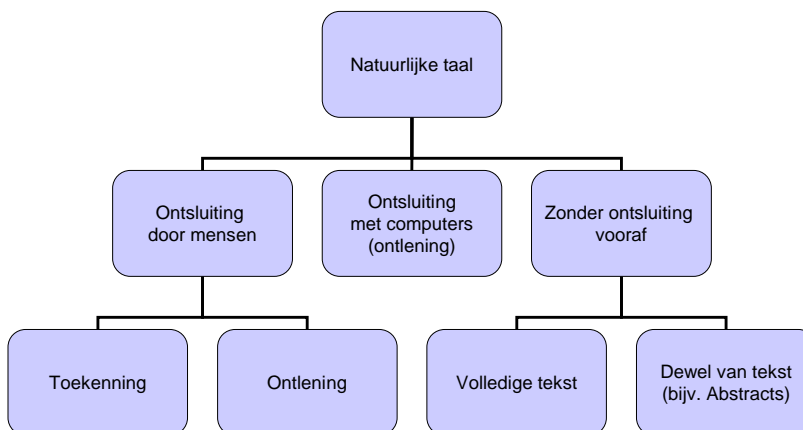


17 feb. 2004

dr. Gerhard J.A. Riesthuis

13

Vrije tekst



17 feb. 2004

dr. Gerhard J.A. Riesthuis

14

Voordelen van vrije tekst

- I Een betere specificiteit in de retrieval
 - Voorbeeld: “geluidshinder door Schiphol in Heiloo”
 - UDC: 628.517.2 : 629.73 (492.62)
[Geluidshinder door luchtvaart in Noord-Holland]



17 feb. 2004

dr. Gerhard J.A. Riesthuis

15

Voordelen van controle

- I Reductie van semantische ambiguïteit
 - Door controle van homografen
- I Bevordering van consistentie bij het weergeven van een onderwerp
 - Door controle van synoniemen
- I Het uitvoeren van veelomvattende zoekacties
 - Door de semantische relaties tussen termen



17 feb. 2004

dr. Gerhard J.A. Riesthuis

16

Kosten

- I Bij gecontroleerde vocabulaires
 - bij de input à de indexer lost de problemen op

- I Bij vrije tekst
 - bij de output à de zoeker lost de problemen op



Twee vragen

1. Kan een menselijke zoeker de problemen oplossen?

2. Kan een zoekprogramma de problemen oplossen?

Maar eerst over de problemen ...



Het probleem van homografen

- bot (been)
 - bot (vis)
 - bot (stomp)
 - bot (brutaal)
- | Tamelijk triviaal probleem in relatief kleine, specialistische domeinen
 - | Minder triviaal in een universeel domein zoals het Internet



17 feb. 2004

dr. Gerhard J.A. Riesthuis

19

Het probleem van synoniemen

- | Echte synoniemen
 - Chemie
 - Scheikunde
- | Ook tamelijk triviaal probleem: gebruik een woordenboek
- | Pseudo-synoniemen: woorden en uitdrukkingen die ongeveer hetzelfde betekenen vormen het echte probleem



17 feb. 2004

dr. Gerhard J.A. Riesthuis

20

Pseudo-synoniemen

- | Komen in vele gedaanten voor
 - Peuter, kleuter, klein kind, vier-jarigen, kinderen uit groep één
 - Toneel, toneelspel, toneelstuk, toneelvoorstelling, toneelspeler
- | Vrije tekst: betere specificiteit
- | Controle: makkelijker “alles” te vinden (betere recall)
- | Ook met vrije-tekstsystemen kunnen goede zoekers redelijke recall bereiken (trunceren, ‘or’)
- | Vereist kennis van het betreffende domein



17 feb. 2004

dr. Gerhard J.A. Riesthuis

21

Het probleem van de veelomvattende zoekacties

- “Regelingen met betrekking tot het houden van huisdieren in Nederlandse gemeenten”
- | Alle gemeenten ...
 - | Alle mogelijke dieren die als huisdier worden gehouden ...
 - | Mogelijke plossing een “explode command”
 - | Vereist echter controle van (pseudo-)synoniemen en semantische relaties



17 feb. 2004

dr. Gerhard J.A. Riesthuis

22

Conclusies

- | Controle speelt vooral een rol voor verbetering van de recall [“het vangen van zoveel mogelijk relevante documenten uit de verzameling relevante documenten”]
- | Vrije-tekstzoeken waardevol voor vragen met een hoge specificiteit



Hoe te controleren?

- | Klassieke informatietalen (met name facetclassificaties en thesauri)
- | Post-controlled vocabulaires



Thesauri

- | Bij het zoeken wordt een synoniem vervangen door de “standaardterm”
- | Relaties meestal beperkt tot
 - Relaties die synoniemen aangeven (UF – USE)
 - Meestal niet tussen samenstellingen en de afzonderlijke termen – Kinderkleding USE Kinderen + Kleding
 - Dit type relaties is belangrijk voor het Nederlands
 - Hierarchische relaties
 - Overige, niet nader gespecificeerde relaties
- | Moet bijgehouden worden



17 feb. 2004

dr. Gerhard J.A. Riesthuis

25

Post-controlled thesauri

- | In principe niet meer dan verzameling van groepen termen met dezelfde of nauw verwante betekenis
- | Vaak ook hierarchische relaties
- | Bij het zoeken wordt de term aangevuld met de andere termen uit de groep waartoe de zoekterm behoort
- | Bijhouden op basis van context in de gevonden documenten



17 feb. 2004

dr. Gerhard J.A. Riesthuis

26

Algemene conclusies

- I Voor goede retrieval nodig:
 - Vrije-tekstzoekmachine
 - Gecontroleerd vocabulaire
 - Kan zowel pre- als post-controlled zijn
 - Bijhouden is een probleem, speciaal in gebieden met snel veranderende vocabulaire
 - Bijhouden met behulp van computers

