

7 Suchen auf XML-Daten

- 7.1 Semistrukturierte Daten in XML
- 7.2 Boolesches Retrieval mit XML-QL
- 7.3 Erweiterung von XML-QL um Keyword-Suche
- 7.4 XXL: Ranked Retrieval auf XML-Daten

Literatur: S. Abiteboul, P. Buneman, D. Suciu,
Data on the Web: From Relations to Semistructured Data
and XML, Morgan Kaufmann, 1999

8. Dezember 2000

Stammvorlesung „Information Retrieval“

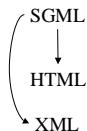
7.1

7.1 Semistrukturierte Daten in XML

XML (Extensible Markup Language) ist **der** (W3C) Standard zur Datenrepräsentation für Datenaustausch, Datenintegration, Datentransformation, E-Commerce (B2B), usw.

Im Gegensatz zu HTML beschreibt XML nur die logische Struktur von Daten; Layout und Präsentation werden durch Style Sheets separat beschrieben (in XSL).

Historische Entwicklung:



XML-zentrierte Technologien:

Definition/Struktur: DTD, XML-Schema
Hyperlinks: Xlink, XPointer
Layout/Filtrn: XSL, CSS
Anfragesprachen:
 XQL, XPATH, XML-QL, XSLT, Quilt, ...
Metadaten: RDF, ...
APIs: DOM, SAX
Spezialisierungen
(z.B. domainspezifische Ontologien):
 MathML, Rosetta, ...

8. Dezember 2000

Stammvorlesung „Information Retrieval“

7.2

Grundkonzepte von XML

- (Frei definierbare) Tags: book, title, author
 - mit Start-Tag: <book> usw.
 - und End-Tag: </book> usw.
 - **Elemente**: <book> ... </book>
Elemente haben einen Namen (book) und einen Inhalt (...)
Elemente können geschachtelt werden
 - Jedes XML-Dokument hat ein Wurzelement und bildet einen Baum
- Elementinhalte können, müssen aber nicht typisiert sein (meist (P)CDATA, also Strings).
Elemente können **Attribute** haben, die jeweils einen Namen und einen Wert (Inhalt) haben, z.B. <article year=1999>. Elemente haben optional Id-Attribute (Element-Ids), auf die innerhalb eines Dokuments mit dem Attribut idref verwiesen werden kann.
Elemente können Hyperlinks als Attribute href haben.

8. Dezember 2000

Stammvorlesung „Information Retrieval“

7.3

XML-Beispiel

```

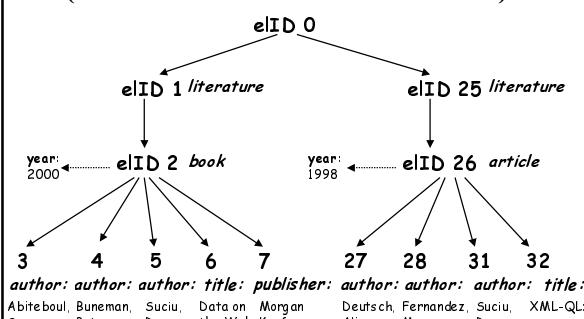
<literature>
  <book year=„2000“>
    <author>Abiteboul, Serge</author>
    <author>Buneman, Peter</author>
    <author>Suciu, Dan</author>
    <title>Data on the Web</title>
    <publisher>Morgan Kaufmann</publisher>
  </book>
  ...
  <article year=„1998“>
    <author>Deutsch, Alin</author>
    <author>Fernandez, Mary</author>
    ...
    <author>Suciu, Dan</author>
    <title>XML-QL: A Query Language for XML</title>
  </article>
  ...
</literature>
  
```

8. Dezember 2000

Stammvorlesung „Information Retrieval“

7.4

XML-Daten als markierte Graphen (mit Elementnamen als Knoten-Labels)



8. Dezember 2000

Stammvorlesung „Information Retrieval“

7.5

Beispiel für Referenzen in XML

```

<person id=„o555“>
  <name>Jane</name>
  <children idref=„o111 o222“/>
</person>
<person id=„o666“>
  <name>Tarzan</name>
  <children idref=„o111 o222“/>
  <homepage href=„http://www.tarzan.biz/home.htm“></homepage>
  <bizportal href=„http://www.tarzan.biz/portal.xml“/>
</person>
<person id=„o111“>
  <name>BillyBoy</name>
  <mother idref=„o555“/><father idref=„o666“/>
</person>
<person id=„o111“>
  <name>Barbie</name>
  <mother idref=„o555“/><father idref=„o666“/>
</person>
  
```

8. Dezember 2000

Stammvorlesung „Information Retrieval“

7.6

DTDs zur XML-Strukturierung

Ein XML-Dokument heißtet **wohlgeformt**, wenn seine Tagstruktur einer korrekten Klammerung entspricht.
Eine DTD (Document Type Definition) ist eine kontextfreie Grammatik zur Beschreibung erlaubter Tagstrukturen (und damit zur Strukturierung von XML-Dokumenten eines Typs). Ein XML-Dokument heißtet **gültig** bzgl. einer DTD wenn seine Tagstruktur durch die DTD generierbar ist.

Beispiel einer DTD:

```
<!DOCTYPE book [
  <!ELEMENT book (title, author*, publisher?, section+)>
  <!ATTLIST book year CDATA #IMPLIED>
  <!ELEMENT title (#PCDATA)>
  <!ELEMENT author (#PCDATA)>
  <!ELEMENT section (#PCDATA | title | section)*>      ]>
```

Weitergehende Schematisierung/Typisierung möglich mit XML-Schema

8. Dezember 2000

Stammvorlesung „Information Retrieval“

7.7

7.2 Boolesches Retrieval mit XML-QL

XML-QL ist eine (von mehreren) Anfragesprache(n) für XML-Daten, die zur Diskussion für die W3C-Standardisierung steht.
Kombination von

- **logischen Bedingungen** über Element- und Attributinhalte (Prädikatenlogik 1. Ordnung a la SQL)
- **regulären Ausdrücken** zum Pattern-Matching von Elementnamen entlang von Pfaden im XML-Datengraph
- + Gruppierung, Aggregation, Strukturtransformation, usw.

Im Gegensatz zu Datenbanksprachen wie SQL bezieht man sich nicht notwendigerweise auf ein festes Strukturschema bzw. muß dieses nicht kennen.

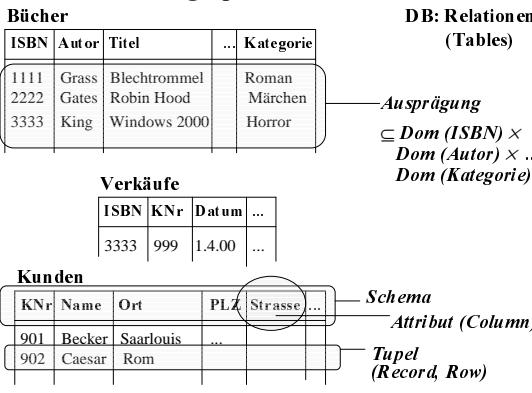
Ein **Anfrageergebnis** ist eine Menge von sich qualifizierenden Pfaden oder Teilgraphen des XML-Datengraphen bzw. daraus konstruierten XML-Dokumenten.

8. Dezember 2000

Stammvorlesung „Information Retrieval“

7.8

Exkurs: Anfragesprache SQL auf Relationen



8. Dezember 2000

Stammvorlesung „Information Retrieval“

7.9

Konzepte von SQL (Structured Query Language)

Anweisungen zum Suchen in Tabellen,
Ändern von Tabellen,
Anlegen neuer Tabellen,
Integritäts sicherung, Zugriffskontrolle, etc.
„stand-alone“ und als API (z.B. in ODBC oder JDBC)

Grobsyntax von SQL-Abfragen (Queries):

```
SELECT (column | expression) { , (column | expression)}
FROM table [correlation_var] {, table [correlation_var]}
[ WHERE search_condition ]
[ GROUP BY column {, column} ]
[ HAVING search_condition ]
```

8. Dezember 2000

Stammvorlesung „Information Retrieval“

7.10

Beispiele für SQL-Abfragen

1) Kunden aus Saarbrücken, die Bücher der Kategorie Politik gekauft haben

```
SELECT K.KNr, K.Name
FROM Kunden K, Verkäufe V, Bücher B
WHERE K.Ort = „Saarbrücken“
AND K.KNr = V.KNr AND V.ISBN = B.ISBN
AND B.Kategorie = „Politik“
```

2) Kunden aus Saarbrücken, die noch nie ein Buch der Kategorie Humor gekauft haben

```
SELECT K.KNr, K.Name
FROM Kunden K
WHERE K.Ort = „Saarbrücken“
AND NOT EXISTS (
  SELECT B.ISBN FROM Bücher B, Verkäufe V
  WHERE B.Kategorie = „Humor“
  AND B.ISBN = V.ISBN AND V.KNr = K.KNr )
```

8. Dezember 2000

Stammvorlesung „Information Retrieval“

7.11

Semantik von SQL-Abfragen

$\mu: \text{SQL} \rightarrow \text{TRK}$ $\text{TRK} = \text{Tupelrelationenkalkül}$
 $\approx \text{Prädikatenlogik 1. Ordnung}$

mit

- $\mu[\text{SELECT } A, B, \dots \text{ FROM } R, S, \dots, T, U, \dots \text{ WHERE } F]$
(wobei A, B, ... Attribute von R, S, ..., nicht aber von T, U, ...)
 $= \{ r.A, s.B, \dots \mid r \in R \wedge s \in S \wedge \dots \wedge \exists t \exists u \dots (t \in T \wedge u \in U \wedge \dots \wedge \mu[F]) \}$
- $\mu^*[F \text{ AND } G] = \mu^*[F] \wedge \mu^*[G]$ • $\mu^*[F \text{ OR } G] = \mu^*[F] \vee \mu^*[G]$
- $\mu^*[\text{NOT } F] = \neg \mu^*[F]$ • $\mu^*[A \theta B] = x.A \theta y.B$
- $\mu^*[\text{EXISTS } S]$ (wobei S: $\text{SELECT } D \text{ FROM } P, Q, \dots \text{ WHERE } H$)
 $= \exists p \exists q \dots (p \in P \wedge q \in Q \wedge \dots \wedge \mu^*[H])$

8. Dezember 2000

Stammvorlesung „Information Retrieval“

7.12

Beispiel für SQL-Semantik

```

SELECT KNr, Name
FROM Kunden K
WHERE K.Ort = „Saarbrücken“
AND NOT EXISTS (SELECT B.ISBN FROM Bücher B, Verkäufe V
                 WHERE B.Kategorie = „Humor“
                   AND B.ISBN = V.ISBN AND V.KNr = K.KNr)

```

$\Downarrow \mu$

$$\begin{aligned}
& \{k.KNr, k.Name \mid k \in \text{Kunden} \wedge \mu^*[K.Ort = \dots \text{ AND NOT EXISTS } \dots] \} \\
&= \{k.KNr, k.Name \mid k \in \text{Kunden} \wedge \mu^*[K.Ort = \dots] \wedge \neg \mu^*[\text{EXISTS } \dots] \} \\
&= \{k.KNr, k.Name \mid k \in \text{Kunden} \wedge k.Ort = \dots \wedge \neg \\
&\quad (\exists b \exists v \mid b \in \text{Bücher} \wedge v \in \text{Verkäufe} \\
&\quad \mu^*[B.Kategorie = \dots \text{ AND } \dots] \wedge \dots) \} \\
&= \{k.KNr, k.Name \mid k \in \text{Kunden} \wedge k.Ort = \dots \wedge \neg \\
&\quad (\exists b \exists v \mid b \in \text{Bücher} \wedge v \in \text{Verkäufe} \wedge \\
&\quad B.Kategorie = \dots \wedge b.ISBN = v.ISBN \wedge v.KNr = k.KNr) \}
\end{aligned}$$

8. Dezember 2000

Stammvorlesung „Information Retrieval“

7.13

Gruppierung und Aggregation in SQL

Beispiel: Verkaufszahlen der Buchkategorien

```

SELECT B.Kategorie, SUM(V.Anzahl) FROM Bücher B, Verkäufe V
WHERE B.ISBN = V.ISBN GROUP BY B.Kategorie

```

ISBN	Autor	Titel	...	Kategorie
1111	Grass	Blechtrömmel		Roman
2222	Gates	Robin Hood		Märchen
3333	Mann	Zauberberg		Roman

ISBN	...	Anzahl
3333	...	1
1111	...	1
2222	...	5
3333	...	2

ISBN	Kategorie	Anzahl
3333	Roman	1
1111	Roman	1
2222	Märchen	5
3333	Roman	2

ISBN	Kategorie	Anzahl
2222	Märchen	5
3333	Roman	1
1111	Roman	1
3333	Roman	2

Kategorie	Anzahl
Märchen	5
Roman	4

8. Dezember 2000

Stammvorlesung „Information Retrieval“

7.14

Konzepte von XML-QL

Eine Anfrage besteht aus drei Blöcken:
Where - In – Construct (\approx Where-From-Select)

Pfadausdrücke:

Elementnamen entlang eines Pfades im XML-Datengraphen
z.B.: <article.author.name>
mit Wildcards bzw. regulären Ausdrücken
z.B. <article.*.name> oder <article*.author>

Elementvariable:

werden an den Elementnamen und -inhalt
am Ende eines Pfades gebunden
z.B.: <author.name> \$n </>
oder an das gesamte Element (mit Unterelementen)
z.B.: <author.name> </> ELEMENT_AS \$n
und können dann mehrfach verwendet werden:
• zur Spezifikation von Filterbedingungen
z.B. \$n Like „%Suci%“ And ...
• zur Konstruktion der Resultatstruktur

8. Dezember 2000

Stammvorlesung „Information Retrieval“

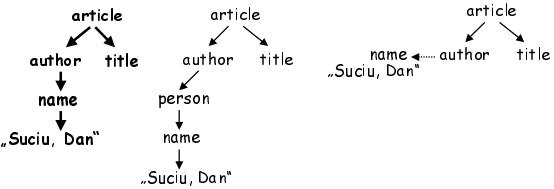
7.15

XML-QL-Beispiel (1)

```

//suche nach Web-bezogenen Artikeln von Dan Suciu aus dem Jahr 1998
WHERE <article year=„1998“>
      <author.name>$n</>
      <title>$t</>
    </>
    IN „www.books/bib.xml“, $n LIKE „Suci%“, $t LIKE „%web%“
CONSTRUCT <result>$n $t</>

```



Treffer

8. Dezember 2000

kein Treffer

kein Treffer

7.16

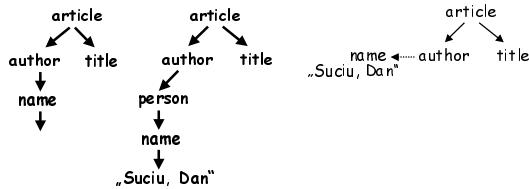
XML-QL-Beispiel (2)

//suche nach Web-bezogenen Artikeln von Dan Suciu aus dem Jahr 1998,
wobei „name“ ein Subelement in beliebiger Tiefe sein kann

```

WHERE <article year=„1998“>
      <author.*.name>$n</>
      <title>$t</>
    </>
    IN „www.books/bib.xml“, $n LIKE „Suci%“, $t LIKE „%web%“
CONSTRUCT <result>$n $t</>

```



Treffer

8. Dezember 2000

Treffer

Stammvorlesung „Information Retrieval“

kein Treffer

7.17

XML-QL: Elementvariable

```

WHERE <$p>
      <title>$t</>
      <$e>Suci, Dan</>
    </>
    IN „www.books/literature.xml“, $p IN {article, book}
CONSTRUCT <$p>
      <title>$t</>
      <$e>Suci, Dan</>
    </>

```

8. Dezember 2000

Stammvorlesung „Information Retrieval“

7.18

XML-QL: Pfadausdrücke

Fragment eines XML-Dokuments:

```
<part>
  <part>
    <part>
      <name>Escort</>
      <brand>Ford</>
    </>
  </>
</>
```

(a) WHERE <part*>
 <name>\$n</>
 <brand>Ford</>
 </>
 IN „www.books/autos.xml“
 CONSTRUCT <result>
 <name>\$n</>
 </>

(b) WHERE <part+.(<name|brand>)>\$n</>
 IN „www.books/autos.xml“
 CONSTRUCT <result>
 <name>\$n</>
 </>

8. Dezember 2000

Stammvorlesung „Information Retrieval“

7.19

XML-QL: Anfragen über mehrere Dokumente (Joins)

// Artikel von Autoren, die ein Buch nach 1995 geschrieben haben

```
WHERE <article>
  <author>
    <firstname>$f</> <lastname>$l</>
  </></> CONTENT_AS $o
  IN „www.a.b.c/bib.xml“,
  <book year=$y>
    <author>
      <firstname>$f</> <lastname>$l</>
    </></>
  IN „www.x.y.z/bib.xml“,
  y > 1995
  CONSTRUCT <article> $a </>
```

8. Dezember 2000

Stammvorlesung „Information Retrieval“

7.20

XML-QL: weitere Features

// suche nach „name“ als Attribut, in beliebiger Tiefe

```
WHERE <article year=„1998“>
  <author* name=$n></>
  <title>$t</>
</>
IN „www.books/bib.xml“, $n LIKE „Suciu%“, $t LIKE „%web%“
CONSTRUCT <result>$n $t</>
```

+ Gruppieren mit geschachtelten Queries
+ Join über Elementnamen
...

8. Dezember 2000

Stammvorlesung „Information Retrieval“

7.21

Speicherung XML-Daten in relationaler DB

Ein Ansatz (von mehreren):

Für jeden Element- oder Attributtyp eine ternäre bzw. binäre Relation, die die von Elementen des entsprechenden Typs ausgehenden Kanten enthält sowie die Element- bzw. Attributwerte (je ein Tripel pro Kante unter Verwendung von ElementIds):

type (source, target, value)

Beispiel: literature

source	target	value
1	2	-
25	26	-

author

source	target	value
3	-	Abiteboul
4	-	Buneman
5	-	Suciu
...		

8. Dezember 2000

Stammvorlesung „Information Retrieval“

7.22

Abbildung von XML-QL auf SQL

// Query in XML-QL:

```
WHERE <article year=„1998“>
  <author name=$n></>
  <title>$t</>
</>
IN „www.books/bib.xml“, $n LIKE „Suciu%“, $t LIKE „%web%“
CONSTRUCT <result>$n $t</>
```

// zugehörige SQL-Query:

```
SELECT B.value, T.value //Name und Titel (Rückgabe)
FROM article A, author B, name N, title T, year Y
WHERE A.target = B.source //Kante article-author
AND A.target = T.source //Kante article-title
AND A.target = Y.source //Kante article-year
AND B.target = N.source //Kante author-name
AND N.value LIKE „Suciu%“
AND Y.value = „1998“ AND T.value LIKE „%web%“
```

8. Dezember 2000

Stammvorlesung „Information Retrieval“

7.23

7.3 Erweiterung von XML-QL um Keyword-Suche

Bisher in XML-QL:

- weniger Wissen über die Struktur macht gezieltes Abfragen von Informationen schwieriger
- Query wird komplexer durch die verschiedenen Fälle: Elementname, Elementinhalt, Attributname, Attributwert

→ Erweiterung:

- Keyword-Suche über Element- und Attributnamen
- Keyword-Suche über (unstrukturierte) Elementinhalt

8. Dezember 2000

Stammvorlesung „Information Retrieval“

7.24

XML-QL-Erweiterung: Contains-Operator	
<pre>CONTAINS(element_var, keyword, depth, expr) //suche nach „name“ als Subelement oder Attribut, bis Tiefe 3 WHERE <article year=„1998“> <author></> ELEMENT_AS \$n <title>\$t</> </> IN „www.books/bib.xml“, CONTAINS (\$n, „Suciu“, 3, any), \$t LIKE „%web%“ CONSTRUCT <result>\$n \$t</> any = el_name OR attr_name OR el_content OR attr_value</pre>	<p>erweitertes invertiertes File:</p> <pre>keywordindex (keyword, elID, depth, location) // location: element, attribut, content, value <literature>, e ID 1, 0, tag> <book>, e ID 2, 1, tag> <year>, e ID 2, 1, attr> ... <author>, e ID 5, 2, tag> <Suciu>, e ID 5, 2, content> ...</pre> <p>Abbildung von Anfragen auf SQL unter Verwendung von „keywordindex“ und (materialisierter) Views</p>

Abbildung Keyword-Suche auf SQL	
<pre>//Query XML-QL mit Stichwort-Suche: WHERE <article year=„1998“> <author></> ELEMENT_AS \$n <title>\$t</> IN „www.books/bib.xml“, CONTAINS (\$n, „Suciu“, 3, any), \$t LIKE „%web%“ CONSTRUCT <result>\$n \$t</> // zugehörige SQL-Query: SELECT B.value, T.value // Name und Titel (Rückgabe) FROM article A, author B, title T, year Y, keywordindex S WHERE S.keyword = „Suciu“ //Keyword-Index AND S.depth < 4 //Tiefebeschränkung AND S.elID = B.source //Kante keywordindex-author AND B.source = A.target //Kante author-article AND A.target = T.source //Kante article-title AND A.target = Y.source //Kante article-year AND Y.value = „1998“ AND T.value LIKE „%web%“ UNION SELECT B.value, T.value FROM ... WHERE ... S.elID = T.source ... UNION ...</pre>	<p>7.27</p> <p>7.28</p> <p>7.29</p> <p>7.30</p> <p>7.31</p> <p>7.32</p> <p>7.33</p> <p>7.34</p> <p>7.35</p> <p>7.36</p> <p>7.37</p> <p>7.38</p> <p>7.39</p> <p>7.40</p> <p>7.41</p> <p>7.42</p> <p>7.43</p> <p>7.44</p> <p>7.45</p> <p>7.46</p> <p>7.47</p> <p>7.48</p> <p>7.49</p> <p>7.50</p> <p>7.51</p> <p>7.52</p> <p>7.53</p> <p>7.54</p> <p>7.55</p> <p>7.56</p> <p>7.57</p> <p>7.58</p> <p>7.59</p> <p>7.60</p> <p>7.61</p> <p>7.62</p> <p>7.63</p> <p>7.64</p> <p>7.65</p> <p>7.66</p> <p>7.67</p> <p>7.68</p> <p>7.69</p> <p>7.70</p> <p>7.71</p> <p>7.72</p> <p>7.73</p> <p>7.74</p> <p>7.75</p> <p>7.76</p> <p>7.77</p> <p>7.78</p> <p>7.79</p> <p>7.80</p> <p>7.81</p> <p>7.82</p> <p>7.83</p> <p>7.84</p> <p>7.85</p> <p>7.86</p> <p>7.87</p> <p>7.88</p> <p>7.89</p> <p>7.90</p> <p>7.91</p> <p>7.92</p> <p>7.93</p> <p>7.94</p> <p>7.95</p> <p>7.96</p> <p>7.97</p> <p>7.98</p> <p>7.99</p> <p>7.100</p> <p>7.101</p> <p>7.102</p> <p>7.103</p> <p>7.104</p> <p>7.105</p> <p>7.106</p> <p>7.107</p> <p>7.108</p> <p>7.109</p> <p>7.110</p> <p>7.111</p> <p>7.112</p> <p>7.113</p> <p>7.114</p> <p>7.115</p> <p>7.116</p> <p>7.117</p> <p>7.118</p> <p>7.119</p> <p>7.120</p> <p>7.121</p> <p>7.122</p> <p>7.123</p> <p>7.124</p> <p>7.125</p> <p>7.126</p> <p>7.127</p> <p>7.128</p> <p>7.129</p> <p>7.130</p> <p>7.131</p> <p>7.132</p> <p>7.133</p> <p>7.134</p> <p>7.135</p> <p>7.136</p> <p>7.137</p> <p>7.138</p> <p>7.139</p> <p>7.140</p> <p>7.141</p> <p>7.142</p> <p>7.143</p> <p>7.144</p> <p>7.145</p> <p>7.146</p> <p>7.147</p> <p>7.148</p> <p>7.149</p> <p>7.150</p> <p>7.151</p> <p>7.152</p> <p>7.153</p> <p>7.154</p> <p>7.155</p> <p>7.156</p> <p>7.157</p> <p>7.158</p> <p>7.159</p> <p>7.160</p> <p>7.161</p> <p>7.162</p> <p>7.163</p> <p>7.164</p> <p>7.165</p> <p>7.166</p> <p>7.167</p> <p>7.168</p> <p>7.169</p> <p>7.170</p> <p>7.171</p> <p>7.172</p> <p>7.173</p> <p>7.174</p> <p>7.175</p> <p>7.176</p> <p>7.177</p> <p>7.178</p> <p>7.179</p> <p>7.180</p> <p>7.181</p> <p>7.182</p> <p>7.183</p> <p>7.184</p> <p>7.185</p> <p>7.186</p> <p>7.187</p> <p>7.188</p> <p>7.189</p> <p>7.190</p> <p>7.191</p> <p>7.192</p> <p>7.193</p> <p>7.194</p> <p>7.195</p> <p>7.196</p> <p>7.197</p> <p>7.198</p> <p>7.199</p> <p>7.200</p> <p>7.201</p> <p>7.202</p> <p>7.203</p> <p>7.204</p> <p>7.205</p> <p>7.206</p> <p>7.207</p> <p>7.208</p> <p>7.209</p> <p>7.210</p> <p>7.211</p> <p>7.212</p> <p>7.213</p> <p>7.214</p> <p>7.215</p> <p>7.216</p> <p>7.217</p> <p>7.218</p> <p>7.219</p> <p>7.220</p> <p>7.221</p> <p>7.222</p> <p>7.223</p> <p>7.224</p> <p>7.225</p> <p>7.226</p> <p>7.227</p> <p>7.228</p> <p>7.229</p> <p>7.230</p> <p>7.231</p> <p>7.232</p> <p>7.233</p> <p>7.234</p> <p>7.235</p> <p>7.236</p> <p>7.237</p> <p>7.238</p> <p>7.239</p> <p>7.240</p> <p>7.241</p> <p>7.242</p> <p>7.243</p> <p>7.244</p> <p>7.245</p> <p>7.246</p> <p>7.247</p> <p>7.248</p> <p>7.249</p> <p>7.250</p> <p>7.251</p> <p>7.252</p> <p>7.253</p> <p>7.254</p> <p>7.255</p> <p>7.256</p> <p>7.257</p> <p>7.258</p> <p>7.259</p> <p>7.260</p> <p>7.261</p> <p>7.262</p> <p>7.263</p> <p>7.264</p> <p>7.265</p> <p>7.266</p> <p>7.267</p> <p>7.268</p> <p>7.269</p> <p>7.270</p> <p>7.271</p> <p>7.272</p> <p>7.273</p> <p>7.274</p> <p>7.275</p> <p>7.276</p> <p>7.277</p> <p>7.278</p> <p>7.279</p> <p>7.280</p> <p>7.281</p> <p>7.282</p> <p>7.283</p> <p>7.284</p> <p>7.285</p> <p>7.286</p> <p>7.287</p> <p>7.288</p> <p>7.289</p> <p>7.290</p> <p>7.291</p> <p>7.292</p> <p>7.293</p> <p>7.294</p> <p>7.295</p> <p>7.296</p> <p>7.297</p> <p>7.298</p> <p>7.299</p> <p>7.300</p> <p>7.301</p> <p>7.302</p> <p>7.303</p> <p>7.304</p> <p>7.305</p> <p>7.306</p> <p>7.307</p> <p>7.308</p> <p>7.309</p> <p>7.310</p> <p>7.311</p> <p>7.312</p> <p>7.313</p> <p>7.314</p> <p>7.315</p> <p>7.316</p> <p>7.317</p> <p>7.318</p> <p>7.319</p> <p>7.320</p> <p>7.321</p> <p>7.322</p> <p>7.323</p> <p>7.324</p> <p>7.325</p> <p>7.326</p> <p>7.327</p> <p>7.328</p> <p>7.329</p> <p>7.330</p> <p>7.331</p> <p>7.332</p> <p>7.333</p> <p>7.334</p> <p>7.335</p> <p>7.336</p> <p>7.337</p> <p>7.338</p> <p>7.339</p> <p>7.340</p> <p>7.341</p> <p>7.342</p> <p>7.343</p> <p>7.344</p> <p>7.345</p> <p>7.346</p> <p>7.347</p> <p>7.348</p> <p>7.349</p> <p>7.350</p> <p>7.351</p> <p>7.352</p> <p>7.353</p> <p>7.354</p> <p>7.355</p> <p>7.356</p> <p>7.357</p> <p>7.358</p> <p>7.359</p> <p>7.360</p> <p>7.361</p> <p>7.362</p> <p>7.363</p> <p>7.364</p> <p>7.365</p> <p>7.366</p> <p>7.367</p> <p>7.368</p> <p>7.369</p> <p>7.370</p> <p>7.371</p> <p>7.372</p> <p>7.373</p> <p>7.374</p> <p>7.375</p> <p>7.376</p> <p>7.377</p> <p>7.378</p> <p>7.379</p> <p>7.380</p> <p>7.381</p> <p>7.382</p> <p>7.383</p> <p>7.384</p> <p>7.385</p> <p>7.386</p> <p>7.387</p> <p>7.388</p> <p>7.389</p> <p>7.390</p> <p>7.391</p> <p>7.392</p> <p>7.393</p> <p>7.394</p> <p>7.395</p> <p>7.396</p> <p>7.397</p> <p>7.398</p> <p>7.399</p> <p>7.400</p> <p>7.401</p> <p>7.402</p> <p>7.403</p> <p>7.404</p> <p>7.405</p> <p>7.406</p> <p>7.407</p> <p>7.408</p> <p>7.409</p> <p>7.410</p> <p>7.411</p> <p>7.412</p> <p>7.413</p> <p>7.414</p> <p>7.415</p> <p>7.416</p> <p>7.417</p> <p>7.418</p> <p>7.419</p> <p>7.420</p> <p>7.421</p> <p>7.422</p> <p>7.423</p> <p>7.424</p> <p>7.425</p> <p>7.426</p> <p>7.427</p> <p>7.428</p> <p>7.429</p> <p>7.430</p> <p>7.431</p> <p>7.432</p> <p>7.433</p> <p>7.434</p> <p>7.435</p> <p>7.436</p> <p>7.437</p> <p>7.438</p> <p>7.439</p> <p>7.440</p> <p>7.441</p> <p>7.442</p> <p>7.443</p> <p>7.444</p> <p>7.445</p> <p>7.446</p> <p>7.447</p> <p>7.448</p> <p>7.449</p> <p>7.450</p> <p>7.451</p> <p>7.452</p> <p>7.453</p> <p>7.454</p> <p>7.455</p> <p>7.456</p> <p>7.457</p> <p>7.458</p> <p>7.459</p> <p>7.460</p> <p>7.461</p> <p>7.462</p> <p>7.463</p> <p>7.464</p> <p>7.465</p> <p>7.466</p> <p>7.467</p> <p>7.468</p> <p>7.469</p> <p>7.470</p> <p>7.471</p> <p>7.472</p> <p>7.473</p> <p>7.474</p> <p>7.475</p> <p>7.476</p> <p>7.477</p> <p>7.478</p> <p>7.479</p> <p>7.480</p> <p>7.481</p> <p>7.482</p> <p>7.483</p> <p>7.484</p> <p>7.485</p> <p>7.486</p> <p>7.487</p> <p>7.488</p> <p>7.489</p> <p>7.490</p> <p>7.491</p> <p>7.492</p> <p>7.493</p> <p>7.494</p> <p>7.495</p> <p>7.496</p> <p>7.497</p> <p>7.498</p> <p>7.499</p> <p>7.500</p> <p>7.501</p> <p>7.502</p> <p>7.503</p> <p>7.504</p> <p>7.505</p> <p>7.506</p> <p>7.507</p> <p>7.508</p> <p>7.509</p> <p>7.510</p> <p>7.511</p> <p>7.512</p> <p>7.513</p> <p>7.514</p> <p>7.515</p> <p>7.516</p> <p>7.517</p> <p>7.518</p> <p>7.519</p> <p>7.520</p> <p>7.521</p> <p>7.522</p> <p>7.523</p> <p>7.524</p> <p>7.525</p> <p>7.526</p> <p>7.527</p> <p>7.528</p> <p>7.529</p> <p>7.530</p> <p>7.531</p> <p>7.532</p> <p>7.533</p> <p>7.534</p> <p>7.535</p> <p>7.536</p> <p>7.537</p> <p>7.538</p> <p>7.539</p> <p>7.540</p> <p>7.541</p> <p>7.542</p> <p>7.543</p> <p>7.544</p> <p>7.545</p> <p>7.546</p> <p>7.547</p> <p>7.548</p> <p>7.549</p> <p>7.550</p> <p>7.551</p> <p>7.552</p> <p>7.553</p> <p>7.554</p> <p>7.555</p> <p>7.556</p> <p>7.557</p> <p>7.558</p> <p>7.559</p> <p>7.560</p> <p>7.561</p> <p>7.562</p> <p>7.563</p> <p>7.564</p> <p>7.565</p> <p>7.566</p> <p>7.567</p> <p>7.568</p> <p>7.569</p> <p>7.570</p> <p>7.571</p> <p>7.572</p> <p>7.573</p> <p>7.574</p> <p>7.575</p> <p>7.576</p> <p>7.577</p> <p>7.578</p> <p>7.579</p> <p>7.580</p> <p>7.581</p> <p>7.582</p> <p>7.583</p> <p>7.584</p> <p>7.585</p> <p>7.586</p> <p>7.587</p> <p>7.588</p> <p>7.589</p> <p>7.590</p> <p>7.591</p> <p>7.592</p> <p>7.593</p> <p>7.594</p> <p>7.595</p> <p>7.596</p> <p>7.597</p> <p>7.598</p> <p>7.599</p> <p>7.600</p> <p>7.601</p> <p>7.602</p> <p>7.603</p> <p>7.604</p> <p>7.605</p> <p>7.606</p> <p>7.607</p> <p>7.608</p> <p>7.609</p> <p>7.610</p> <p>7.611</p> <p>7.612</p> <p>7.613</p> <p>7.614</p> <p>7.615</p> <p>7.616</p> <p>7.617</p> <p>7.618</p> <p>7.619</p> <p>7.620</p> <p>7.621</p> <p>7.622</p> <p>7.623</p> <p>7.624</p> <p>7.625</p> <p>7.626</p> <p>7.627</p> <p>7.628</p> <p>7.629</p> <p>7.630</p> <p>7.631</p> <p>7.632</p> <p>7.633</p> <p>7.634</p> <p>7.635</p> <p>7.636</p> <p>7.637</p> <p>7.638</p> <p>7.639</p> <p>7.640</p> <p>7.641</p> <p>7.642</p> <p>7.643</p> <p>7.644</p> <p>7.645</p> <p>7.646</p> <p>7.647</p> <p>7.648</p> <p>7.649</p> <p>7.650</p> <p>7.651</p> <p>7.652</p> <p>7.653</p> <p>7.654</p> <p>7.655</p> <p>7.656</p> <p>7.657</p> <p>7.658</p> <p>7.659</p> <p>7.660</p> <p>7.661</p> <p>7.662</p> <p>7.663</p> <p>7.664</p> <p>7.665</p> <p>7.666</p> <p>7.667</p> <p>7.668</p> <p>7.669</p> <p>7.670</p> <p>7.671</p> <p>7.672</p> <p>7.673</p> <p>7.674</p> <p>7.675</p> <p>7.676</p> <p>7.677</p> <p>7.678</p> <p>7.679</p> <p>7.680</p> <p>7.681</p> <p>7.682</p> <p>7.683</p> <p>7.684</p> <p>7.685</p> <p>7.686</p> <p>7.687</p> <p>7.688</p> <p>7.689</p> <p>7.690</p> <p>7.691</p> <p>7.692</p> <p>7.693</p> <p>7.694</p> <p>7.695</p> <p>7.696</p> <p>7.697</p> <p>7.698</p> <p>7.699</p> <p>7.700</p> <p>7.701</p> <p>7.702</p> <p>7.703</p> <p>7.704</p> <p>7.705</p> <p>7.706</p> <p>7.707</p> <p>7.708</p> <p>7.709</p> <p>7.710</p> <p>7.711</p> <p>7.712</p> <p>7.713</p> <p>7.714</p> <p>7.715</p> <p>7.716</p> <p>7.717</p> <p>7.718</p> <p>7.719</p> <p>7.720</p> <p>7.721</p> <p>7.722</p> <p>7.723</p> <p>7.724</p> <p>7.725</p> <p>7.726</p> <p>7.727</p> <p>7.728</p> <p>7.729</p> <p>7.730</p> <p>7.731</p> <p>7.732</p> <p>7.733</p> <p>7.734</p> <p>7.735</p> <p>7.736</p> <p>7.737</p> <p>7.738</p> <p>7.739</p> <p>7.740</p> <p>7.741</p> <p>7.742</p> <p>7.743</p> <p>7.744</p> <p>7.745</p> <p>7.746</p> <p>7.747</p> <p>7.748</p> <p>7.749</p> <p>7.750</p> <p>7.751</p> <p>7.752</p> <p>7.753</p> <p>7.754</p> <p>7.755</p> <p>7.756</p> <p>7.757</p> <p>7.758</p> <p>7.759</p> <p>7.760</p> <p>7.761</p> <p>7.762</p> <p>7.763</p> <p>7.764</p> <p>7.765</p> <p>7.766</p> <p>7.767</p> <p>7.768</p> <p>7.769</p> <p>7.770</p> <p>7.771</p> <p>7.772</p> <p>7.773</p> <p>7.774</p> <p>7.775</p> <p>7.776</p> <p>7.777</p> <p>7.778</p> <p>7.779</p> <p>7.780</p> <p>7.781</p> <p>7.782</p> <p>7.783</p> <p>7.784</p> <p>7.785</p> <p>7.786</p> <p>7.787</p> <p>7.788</p> <p>7.789</p> <p>7.790</p> <p>7.791</p> <p>7.792</p> <p>7.793</p> <p>7.794</p> <p>7.795</p> <p>7.796</p> <p>7.797</p> <p>7.798</p> <p>7.799</p> <p>7.800</p> <p>7.801</p> <p>7.802</p> <p>7.803</p> <p>7.804</p> <p>7.805</p> <p>7.806</p> <p>7.807</p> <p>7.808</p> <p>7.809</p> <p>7.810</p> <p>7.811</p> <p>7.812</p> <p>7.813</p> <p>7.814</p> <p>7.815</p> <p>7.816</p> <p>7.817</p> <p>7.818</p> <p>7.819</p> <p>7.820</p> <p>7.821</p> <p>7.822</p> <p>7.823</p> <p>7.824</p> <p>7.825</p> <p>7.826</p> <p>7.827</p> <p>7.828</p> <p>7.829</p> <p>7.830</p> <p>7.831</p> <p>7.832</p> <p>7.833</p> <p>7.834</p> <p>7.835</p> <p>7.836</p> <p>7.837</p> <p>7.838</p> <p>7.839</p> <p>7.840</p> <p>7.841</p> <p>7.842</p> <p>7.843</p> <p>7.844</p> <p>7.845</p> <p>7.846</p> <p>7.847</p> <p>7.848</p> <p>7.849</p> <p>7.850</p> <p>7.851</p> <p>7.852</p> <p>7.853</p> <p>7.854</p> <p>7.855</p> <p>7.856</p> <p>7.857</p> <p>7.858</p> <p>7.859</p> <p>7.860</p> <p>7.861</p> <p>7.862</p> <p>7.863</p> <p>7.864</p> <p>7.865</p> <p>7.866</p> <p>7.867</p> <p>7.868</p> <p>7.869</p> <p>7.870</p> <p>7.871</p> <p>7.872</p> <p>7.873</p> <p>7.874</p> <p>7.875</p> <p>7.876</p> <p>7.877</p> <p>7.878</p> <p>7.879</p> <p>7.880</p> <p>7.881</p> <p>7.882</p> <p>7.883</p> <p>7.884</p> <p>7.885</p> <p>7.886</p> <p>7.887</p> <p>7.888</p> <p>7.889</p> <p>7.890</p> <p>7.891</p> <p>7.892</p> <p>7.893</p> <p>7.894</p> <p>7.895</p> <p>7.896</p> <p>7.897</p> <p>7.898</p> <p>7.899</p> <p>7.900</p> <p>7.901</p> <p>7.902</p> <p>7.903</p> <p>7.904</p> <p>7.905</p> <p>7.906</p> <p>7.907</p> <p>7.908</p> <p>7.909</p> <p>7.910</p> <p>7.911</p> <p>7.912</p> <p>7.913</p> <p>7.914</p> <p>7.915</p> <p>7.916</p> <p>7.917</p> <p>7.918</p> <p>7.919</p> <p>7.920</p> <p>7.921</p> <p>7.922</p> <p>7.923</p> <p>7.924</p> <p>7.925</p> <p>7.926</p> <p>7.927</p> <p>7.928</p> <p>7.929</p> <p>7.930</p> <p>7.931</p> <p>7.932</p> <p>7.933</p> <p>7.934</p> <p>7.935</p> <p>7.936</p> <p>7.937</p> <p>7.938</p> <p>7.939</p> <p>7.940</p> <p>7.941</p> <p>7.942</p> <p>7.943</p> <p>7.944</p> <p>7.945</p> <p>7.946</p> <p>7.947</p> <p>7.948</p> <p>7.949</p> <p>7.950</p> <p>7.951</p> <p>7.952</p> <p>7.953</p> <p>7.954</p> <p>7.955</p> <p>7.956</p> <p>7.957</p> <p>7.958</p> <p>7.959</p> <p>7.960</p> <p>7.961</p> <p>7.962</p> <p>7.963</p> <p>7.964</p> <p>7.965</p> <p>7.966</p> <p>7.967</p> <p>7.968</p> <p>7.969</p> <p>7.970</p> <p>7.971</p> <p>7.972</p> <p>7.973</p> <p>7.974</p> <p>7.975</p> <p>7.976</p> <p>7.977</p> <p>7.978</p> <p>7.979</p> <p>7.980</p> <p>7.981</p> <p>7.982</p> <p>7.983</p> <p>7.984</p> <p>7.985</p> <p>7.986</p> <p>7.987</p> <p>7.988</p> <p>7.989</p> <p>7.990</p> <p>7.991</p> <p>7.992</p> <p>7.993</p> <p>7.994</p> <p>7.995</p> <p>7.996</p> <p>7.997</p> <p>7.998</p> <p>7.999</p> <p>7.100</p> <p>7.101</p> <p>7.102</p> <p>7.103</p> <p>7.104</p> <p>7.105</p> <p>7.106</p> <p>7.107</p> <p>7.108</p> <p>7.109</p> <p>7.110</p> <p>7.111</p> <p>7.112</p> <p>7.113</p> <p>7.114</p> <p>7.115</p> <p>7.116</p> <p>7.117</p> <p>7.118</p> <p>7.119</p> <p>7.120</p> <p>7.121</p> <p>7.122</p> <p>7.123</p> <p>7.124</p> <p>7.125</p> <p>7.126</p> <p>7.127</p> <p>7.128</p> <p>7.129</p> <p>7.130</p> <p>7.131</p> <p>7.132</p> <p>7.133</p> <p>7.134</p> <p>7.135</p> <p>7.136</p> <p>7.137</p> <p>7.138</p> <p>7.139</p> <p>7.140</p> <p>7.141</p> <p>7.142</p> <p>7.143</p> <p>7.144</p> <p>7.145</p> <p>7.146</p> <p>7.147</p> <p>7.148</p> <p>7.149</p> <p>7.150</p> <p>7.151</p> <p>7.152</p> <p>7.1</p>



8. Dezember 2000

Stammvorlesung „Information Retrieval“

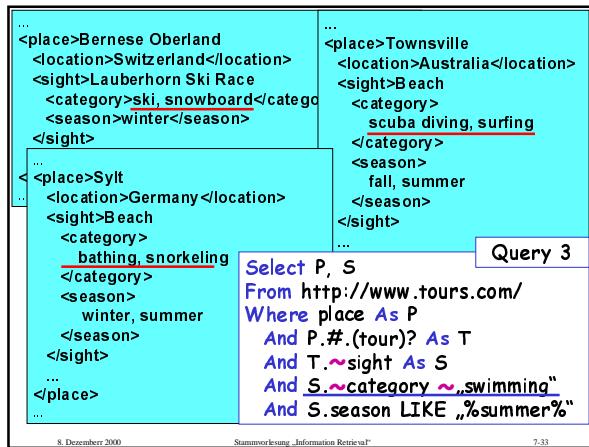
7.31



8. Dezember 2000

Stammvorlesung „Information Retrieval“

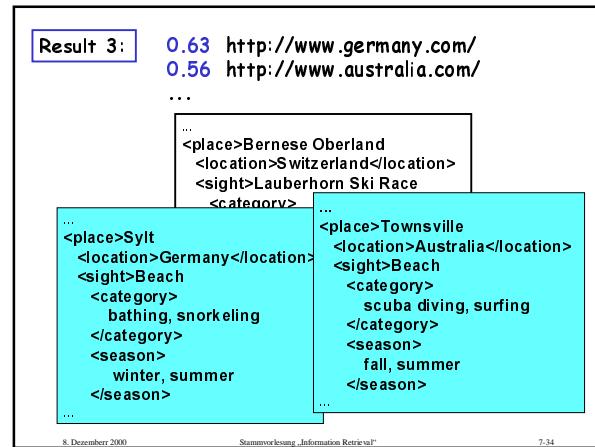
7.32



8. Dezember 2000

Stammvorlesung „Information Retrieval“

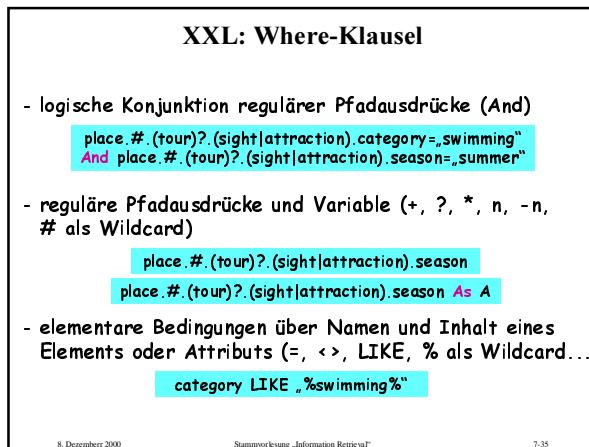
7.33



8. Dezember 2000

Stammvorlesung „Information Retrieval“

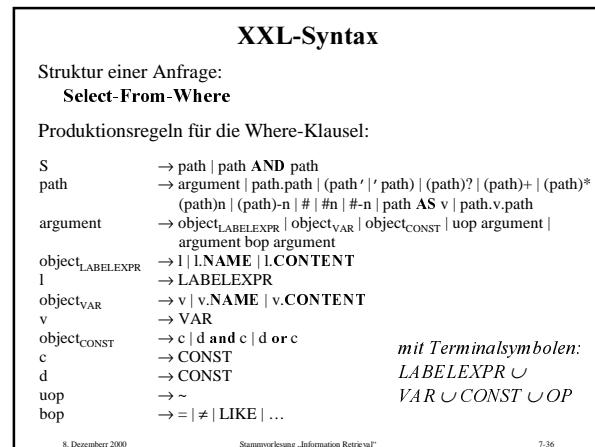
7.34



8. Dezember 2000

Stammvorlesung „Information Retrieval“

7.35



8. Dezember 2000

Stammvorlesung „Information Retrieval“

7.36

Vereinfachte XXL-Semantik (1)

Ein **Element** ist ein 5-Tupel der Form (oid, label, content, elements, attributes) mit einem Bezeichner oid, einem String label, einem String content, einer Liste elements von weiteren Elementen, die auch leer sein kann, und einer Liste attributes von Attributen.
 Ein **Dokument** ist ein Element, das in keinem anderen Element enthalten ist. Für ein Element x bezeichnen x.oid, x.label, x.content, x.elements die entsprechenden Komponenten; x.elements[i] bezeichnet das i-te Element.
 Ein **Datengraph** ist der einem Dokument x entsprechende Graph $G=(V,E)$ mit Elementen als Knotenmenge V und einer Kantenmenge E, die eine Kante von a nach b enthält, wenn b in a.elements ist. Knoten sind mit den Labels der Elemente markiert. Die Reihenfolge in a.elements wird im Graphen nicht mehr berücksichtigt. Ein Pfad p im Graph entspricht einer Label-Sequenz.

Notation:

- LABEL $\subseteq \Sigma^*$ die Menge der in x vorkommenden Labels,
- PATH $\subseteq V^+$ die Menge der Pfade in G
(der Länge k für k Knoten und k-1 Kanten für k=1, 2, ...),
- VAR eine Menge von Knotenvariablen.

8. Dezember 2000

Stammvorlesung „Information Retrieval“

7.37

Vereinfachte XXL-Semantik (2)

LABELEXPR ist die Menge $(\Sigma \cup \{\%\})^*$ mit dem Wildcard-Zeichen %.
 PATHEXPR, ist die Menge der regulären Pfadausdrücke, d.h. die kleinste Menge mit den folgenden Eigenschaften:

- Jedes Wort aus LABELEXPR ist ein Pfadausdruck.
- Jede Variable A ∈ VAR ist ein Pfadausdruck (für Pfad der Länge 1).
- Falls S und T Pfadausdrücke sind, dann sind auch ST (Konkatenation), S|T (Vereinigung), (S) (Klammerung), (S)? (optionales Vorkommen), (S)+ (mehrfaches Vorkommen), (S)* (optional mehrfaches Vorkommen) und # (Abk. für (%)*) Pfadausdrücke.

Eine **Variablenbindung** zu einem gegebenen Datengraphen $G=(V,E)$ ist eine Abbildung $v: VAR \rightarrow V$. $v[A]$ bezeichnet den Wert der Variablen A.

8. Dezember 2000

Stammvorlesung „Information Retrieval“

7.38

Vereinfachte XXL-Semantik (3)

Ein Pfad p eines Datengraphen G erfüllt einen Pfadausdruck E bei Variablenbindung v, wenn folgendes gilt:

- Falls E ein Labelausdruck ist, dann muß p Länge 1 haben und dem Stringmuster E $\in (\Sigma \cup \{\%\})^*$ entsprechen.
- Falls E eine Variable A ist, dann muß p Länge 1 haben und $v[A] = \text{first}(p)$ sein.
- Falls E die Form ST hat, dann muß p zerlegbar sein in pre(p) und suf(p) mit $p = \text{pre}(p)suf(p)$, $|\text{pre}(p)| \geq 1$, $|\text{suf}(p)| \geq 1$, so daß pre(p) S erfüllt und suf(p) T erfüllt.
- Falls E die Form S|T hat, dann muß p S oder T erfüllen.
- Falls E die Form (S) hat, dann muß p S erfüllen.
- Falls E die Form (S)? hat, dann muß p leer sein oder S erfüllen.
- Falls E die Form (S)+ hat, dann muß p zerlegbar sein in pre(p) und suf(p) mit $p = \text{pre}(p)suf(p)$, $|\text{pre}(p)| \geq 1$, $|\text{suf}(p)| \geq 1$, so daß pre(p) S erfüllt und suf(p) (S)+ erfüllt.
- Falls E die Form (S)* hat, dann muß p leer sein oder (S)+ erfüllen.
- Falls E die Form # hat, erfüllt p E.

8. Dezember 2000

Stammvorlesung „Information Retrieval“

7.39

Vereinfachte XXL-Semantik (4)

Eine **elementare Query** Q auf einem Datengraphen $G=(V,E)$ ist von der Form 1) E As A oder 2) [E As] A op b oder 3) [E As] A op B mit $E \in \text{PATHEXPR}$, $\text{op} \in \{<, >, \leq, \geq, =, \neq, \text{LIKE}\}$, Variablen A, B und Konst. b. Eine **Query** Q auf $G=(V,E)$ ist eine Konjunktion von elementaren Queries Q_1, \dots, Q_m auf G ist:

Die **Semantik einer elementaren Query** Q auf G bei Variablenbelegung v ist:
 $\mu_v(Q, G) = \{p \in V^+ \mid p \text{ erfüllt } E \text{ und es gilt }$
 $\text{last}(p).\text{content} \text{ op } v[A] \text{ op } b \text{ bzw. last}(p).\text{content} \text{ op } v[A] \text{ op } v[B]\}$
 Die **Semantik einer Query** $Q = Q_1 \& Q_2 \& \dots \& Q_m$ mit elementaren Queries Q_1, \dots, Q_m auf G ist:
 $\mu(Q, G) = \{\text{Teilgraphen } G^i \text{ von } G \mid G^i \text{ ist die Vereinigung von Pfaden } p_1, \dots, p_m \in V^+ \text{ und es gibt eine Variablenbelegung } v, \text{ so daß}$
 $\text{für alle } i=1, \dots, m \text{ gilt: } p_i \in \mu_v(Q_i, G)\}$

8. Dezember 2000

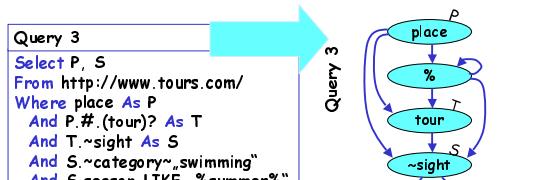
Stammvorlesung „Information Retrieval“

7.40

XXL: Pfadausdruck → NEA

Where-Klausel einer Query

- Ähnlichkeitsoperator \sim



8. Dezember 2000

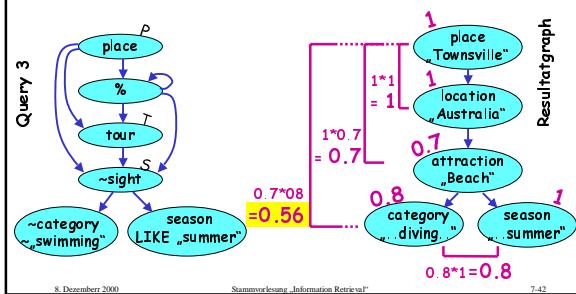
Stammvorlesung „Information Retrieval“

7.41

XXL: Ähnlichkeitsoperator \sim (1)

Where-Klausel einer Query

- Ähnlichkeitsoperator \sim
- Regeln zur Relevanzbewertung von Knoten und Pfaden



8. Dezember 2000

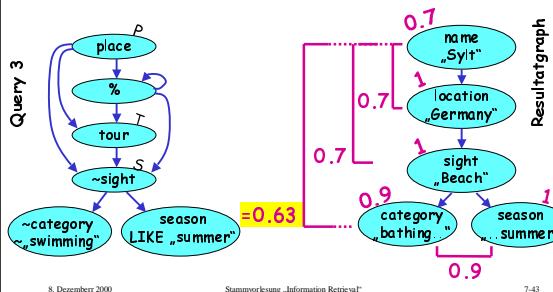
Stammvorlesung „Information Retrieval“

7.42

XXL: Ähnlichkeitsoperator ~ (2)

Where-Klausel einer Query

- Ähnlichkeitsoperator \sim
- Regeln zur Relevanzbewertung von Knoten und Pfaden



XXL: Ähnlichkeitsoperator ~ (3)

Query 3: Where place As P

```
And P.#.(tour)? As T And T.~sight As S
And S.~category="swimming"
And S.season LIKE "%summer%"
```

Result 3: 0.63 http://www.germany.com/
0.56 http://www.australia.com/

```
<place>Sylt
<location>Germany</location>
<sight>Beach
<category>bathing, snorkeling
</category>
<season>winter, summer
</season>
</sight>
...
<place>Bernese Oberland
<location>Switzerland</location>
<sight>Interlaken
<category>skiing, bathing, snorkeling
<season>winter, summer
</season>
</sight>
...
<place>Townsville
<location>Australia</location>
<sight>Beach
<category>scuba diving, surfing
<category>
<season>fall, summer
</season>
</sight>
...
```

8. Dezember 2000 Stammtagskongress „Information Retrieval“ 7.44

XXL-Semantik mit Ähnlichkeitsoperator ~ (1)

Relevanz-Score für primitive Bedingungen:

- Für Elementinhaltsvergleiche basiert der Score auf der tf*idf-Formel.
- Für Element- und Attributnamensvergleiche basiert der Score wahlweise auf der
 - Levenshtein-String-Ähnlichkeit oder
 - einer „semantischen“ Ähnlichkeitsfunktion auf einem Ontologiebaum für Elementnamen
- Für Attributwertevergleiche ist der Score typspezifisch (z.B. basierend auf Intervallen bei Geldbeträgen oder Datumsangaben oder berechneten Distanzen, Fahrzeiten usw. bei Ortsangaben)

8. Dezember 2000 Stammtagskongress „Information Retrieval“ 7.45

XXL-Semantik mit Ähnlichkeitsoperator ~ (2)

Relevanz-Score für Pfade:

- Für Pfade p1 und p2 mit Scores π_1 und π_2 bzgl. Bedingungen c1 und c2 ist der Score von p1.p2 bzgl. "c1.c2" $\pi_1 * \pi_2$.
- Für Pfad p mit Score π bzgl. c ist der Score von p bzgl. "c?" π .
- Für Pfade p1, p2, ..., pm mit Scores $\pi_1, \pi_2, \dots, \pi_m$ bzgl. c ist der Score von p1.p2.pm bzgl. "c+", "c*", "cn" und "c-n" $\pi_1 * \pi_2 * \dots * \pi_m$ bzw.
- $\pi_1 * \pi_2 * \dots * \pi_m$ falls p1.p2.pm nicht leer ist, 1 sonst bzw. $\pi_1 * \pi_2 * \dots * \pi_m$ falls n=m, 0 sonst bzw.
- $\pi_1 * \pi_2 * \dots * \pi_m$ falls m < n, 0 sonst.
- Für Pfad p und Bedingung "%" (z.B. "(%)*" bzw. "#") ist der Score 1.

Relevanz-Score für Teilgraphen:

- Für Teilgraphen g1 und g2 mit Scores π_1 und π_2 bzgl. c1 und c2 ist der Score von $g1 \cup g2$ bzgl. "c1 | c2" $\pi_1 + \pi_2 - \pi_1 * \pi_2$ und der Score von $g1 \cup g2$ bzgl. "c1 & c2" $\pi_1 * \pi_2$.
- Für Teilgraphen g1 und g2 mit Scores π_1 und π_2 bzgl. c1 und c2 ist der Score von $g1.g2$ bzgl. "c1.c2" $\pi_1 * \pi_2$.

8. Dezember 2000 Stammtagskongress „Information Retrieval“ 7.46

XXL-Semantik mit Ähnlichkeitsoperator ~ (3)

Für einen Ontologiebaum (V, E) mit gewichteten Kanten ist die Distanz zwischen zwei Knoten, $dist: V \times V \rightarrow [0..1]$, wie folgt definiert:

- Für Knoten v_1 und v_2 auf demselben Wurzel-Blatt-Pfad ist $dist(v_1, v_2)$ die Anzahl der Kanten zwischen v_1 und v_2 .
- Für Geschwisterknoten v_1, v_2 mit Vater p ist $siblingdist(v_1, v_2) = |weight(p, v_1) - weight(p, v_2)| / |children(p)|$
- Für beliebige Knoten v_1, v_2 mit dem kleinsten gemeinsamen Vorfahren p , der Kinder p_1 und p_2 mit Pfaden $p.p_1 \dots v_1$ und $p.p_2 \dots v_2$ hat $dist(v_1, v_2) = (length(p_1, v_1) + length(p_2, v_2) + siblingdist(p_1, p_2) + 1) / maxdist(v_1, v_2)$, wobei $maxdist(v_1, v_2) = length(r, v_1) + length(r, v_2)$ ein Normierungsfaktor und r die Wurzel des Baums ist.
- Die Ähnlichkeit zweier Knoten ist dann: $sim(v_1, v_2) = 1 - dist(v_1, v_2)$.

8. Dezember 2000 Stammtagskongress „Information Retrieval“ 7.47

Implementierung mit Oracle8i interMedia

Abbildung von XML-Dokumenten auf relationale Tabellen:

Documents (URL, root_oid)	Attributes (oid, name, value)
Elements (oid, name, content)	Edges (oid, parent_oid)

```
Select P, S From http://www.tours.com/
Where place As P And P.#.(tour)? As T And T.~sight As S
And S.~category="swimming" And S.season="summer"
```

Abbildung von XXL-Anfragen auf SQL und Oracle8i interMedia

→ Testen elementarer Bedingungen + Aufbau des Resultatgraphs

```
Select e, SCORE(1)+SCORE(2)-SCORE(1)*SCORE(2)
From Elements e, Edges v
Where CONTAINS(e.name, SYN(category), 1)>0
And CONTAINS(e.content, "swimming", 2)>0
And e.oid = v.oid
And v.parent_oid=a current leaf in resultgraph
```

8. Dezember 2000 Stammtagskongress „Information Retrieval“ 7.48

