

# TAKSONOMI MODEL-MODEL SISTEM TEMU KEMBALI INFORMASI

MATERI PERKULIAHAN : *INFORMATION RETRIEVAL SYSTEM* KE-4

Disusun Oleh :

Nama : Nuning Kurniasih, S.Sos., M.Hum.

NIP. 197606252000122001

Departemen Ilmu Informasi dan Perpustakaan

Fakultas Ilmu Komunikasi

Universitas Padjadjaran

Ditulis Pertama Tahun 2005, Revisi Januari 2014



# JENIS SISTEM TEMU KEMBALI INFORMASI

1. Lokal

2. Global

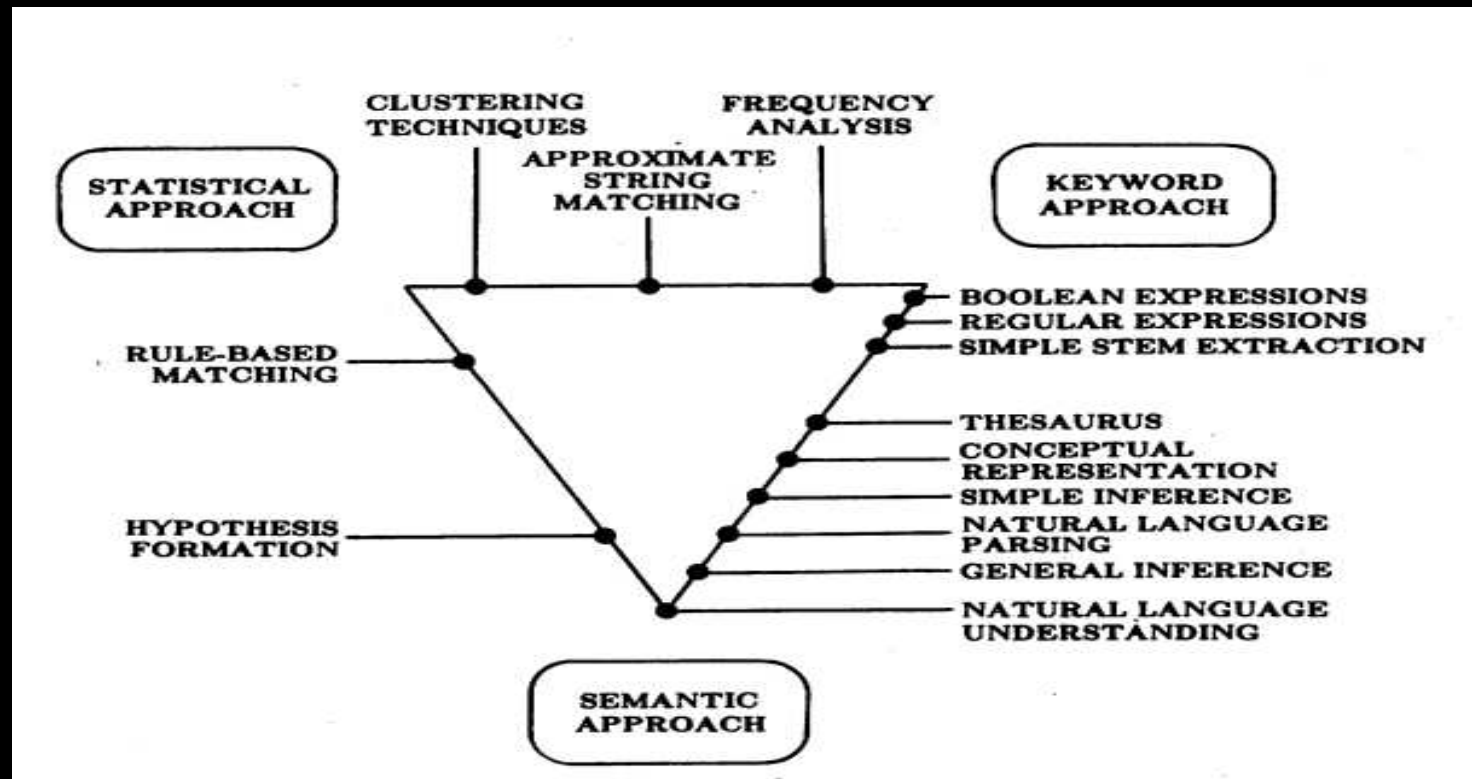
Baca kembali materi perkuliahan ke-2

# World Wide Web



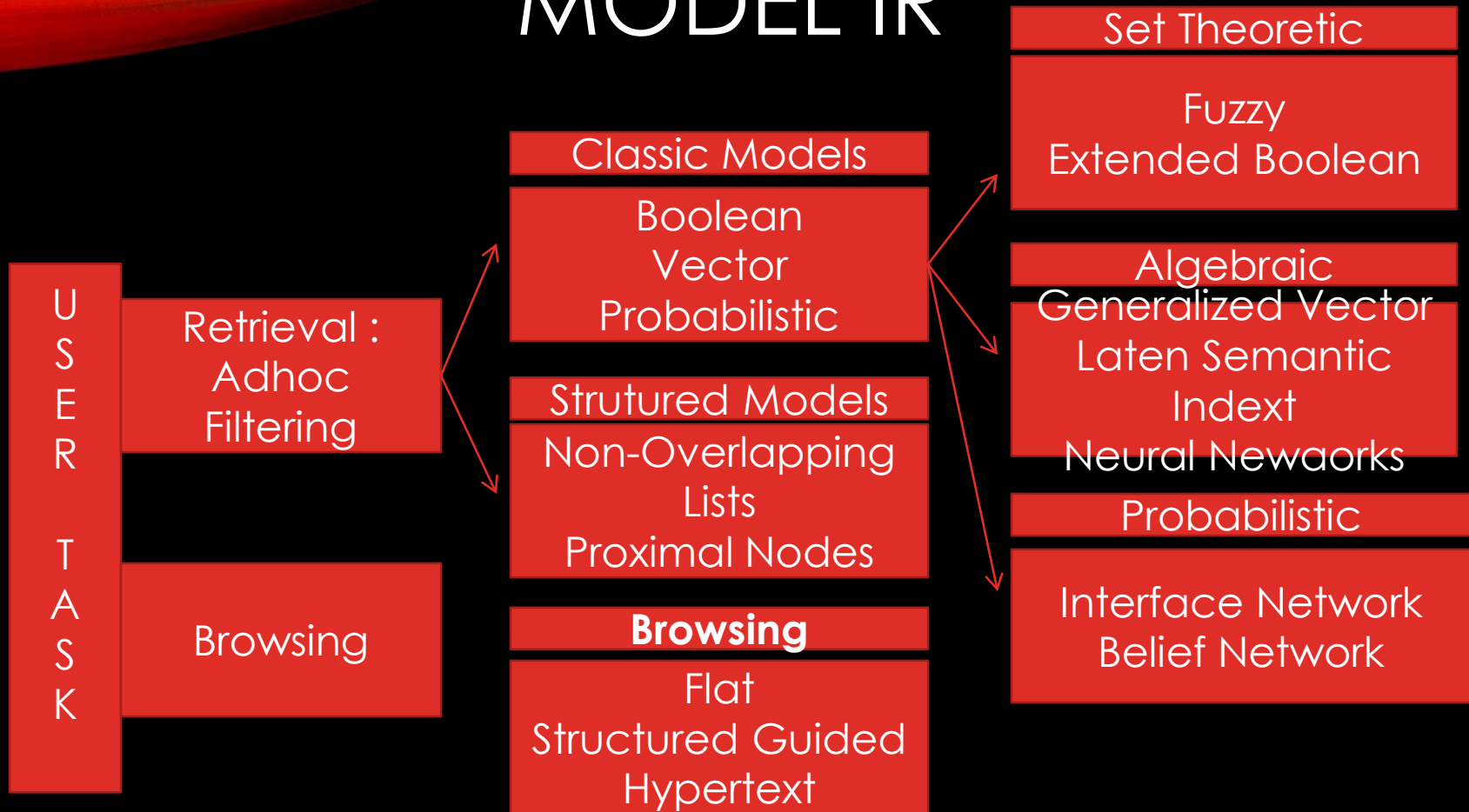
<http://blog.law.cornell.edu>

# THE INFORMATION RETRIEVAL TRIANGLE



Mc Cune, Brian P. 1985.

# MODEL IR



Baeza-Yates & Ribeiro Neto

# LOGICAL VIEW OF DOC AND RETRIEVAL TASK

## LOGICAL VIEW OF DOCUMENTS

U  
S  
E  
R  
  
T  
A  
S  
K

|           | Index Term   | Full Text  | Full Text + Structure            |
|-----------|--|--|----------------------------------|
| Retrieval | Classic<br>Set Theoretic<br>Algebraic<br>Probabilistic | Classic<br>Set Theoretic<br>Algebraic<br>Probabilistic | Structured                       |
| Browsing  | Flat   | Flat<br>Hypertext                                      | Structured<br>Guide<br>Hypertext |

Baeza-Yates & Ribeiro Neto

## KONSEP DASAR MODEL-MODEL IR KLASIK

- ▶ Dokumen direpresentasikan dengan seperangkat *keyword* atau *index term* yang representative.
- ▶ Sebuah item indeks adalah kata-kata dalam dokumen yang mudah diingat untuk tema umum sebuah dokumen.
- ▶ *Index term* biasanya kata benda.
- ▶ *Search engines* mengasumsikan semua kata adalah *index term* (merepresentasikan full text)
- ▶ Tidak semua *term* dapat merepresentasikan isi dokumen, sebagian *term* dapat mengidentifikasi dokumen dengan sempit.



## KONSEP DASAR MODEL-MODEL IR KLASIK

- ▶  $K_j$  adalah *index term*
- ▶  $d_j$  adalah sebuah dokumen
- ▶  $t$  adalah total dari jumlah dokumen
- ▶  $K = (k_1, k_2, \dots, k_t)$  adalah seluruh *index term*
- ▶  $w_{ij} \geq 0$  is penimbang yang berasosiasi dengan  $(k_i, d_j)$
- ▶  $w_{ij} = 0$  mengindikasikan bahwa term tidak berhubungan dengan dokumen
- ▶  $\text{vec}(d_j) = (w_{1j}, w_{2j}, \dots, w_{tj})$  adalah penimbang *vector* yang berasosiasi dengan dokumen  $d_j$ .
- ▶  $g_i(\text{vec}(d_j)) = w_{ij}$  adalah sebuah fungsi dengan pengembalian pertimbangan yang berasosiasi dengan pasangannya associated with pair  $(k_j, d_j)$



## MODEL TEMU KEMBALI INFORMASI KLASIK

- Boolean Query
  - Operator Boolean Logic : AND, OR, NOT
- Vector Query
  - Seberapa mirip dokumen dengan?
    - [(java 3) (compiler 2) (unix 1) (linus 1)]
- Probabilistic Query
  - Probabilistik Relevan -  $Pr(rel)$
  - Probabilistik Non Relevan -  $Pr(non\ rel)$

# MODEL PENELITIAN BOOLEAN

Temu kembali berbasis pada kriteria pengambilan keputusan dua kondisi (binary) tanpa ide atau pencocok yang pasial.

Tidak ada dalam dokumen yang disediakan.

Informasi perlu diterjemahkan kedalam sebuah *Boolean expression*.

Operator *Boolean Logic* diekspresikan dengan

Logical AND

Logical OR

Logical NOT

# BOOLEAN STANDAR

|                | Standard Boolean   |
|----------------|--|
| <b>Goal</b>    | <ul style="list-style-type: none"> <li>• Capture conceptual structure and contextual information</li> </ul>  |
| <b>Methods</b> | <ul style="list-style-type: none"> <li>• Coordination: AND, OR, NOT</li> <li>• Proximity</li> <li>• Fields</li> <li>• Stemming / Truncation</li> </ul>   |
| <b>(+)</b>     | <ul style="list-style-type: none"> <li>• Easy to implement</li> <li>• Computationally efficient<br/>=&gt; all the major on-line databases use it</li> <li>• Expressiveness and Clarity<br/>Synonym specifications (OR-clauses) and phrases (AND-clauses).</li> </ul>   |
| <b>(-)</b>     | <ul style="list-style-type: none"> <li>• Difficult to construct Boolean queries.</li> <li>• All or nothing<br/>AND too severe, and OR does not differentiate enough.</li> <li>• Difficult to control output: Null output &lt;--&gt; Overload.</li> <li>• No ranking</li> <li>• No weighting of index or query terms</li> <li>• No uncertainty measure</li> </ul> |

Ringkasan dari Karakteristik Standar Pendekatan Boolean disertai Kelebihan dan Kelemahannya  
<http://comminfo.rutgers.edu>

# LOGICAL AND

- Memperbolehkan penelusur untuk menggunakan pernyataan query ke dalam dua atau lebih konsep sehingga hasil penelusuran menjadi lebih terbatas.
- Formula pernyataan sederhana A AND B
- Contoh untuk menelusur *marketing and library*, kita memformulasikan pernyataan penelusuran dengan :  
*marketing **AND** library*
- Dengan query tersebut maka kita akan menemukan dokumen yang mengandung unsur marketing dan perpustakaan saja, dan tidak untuk mendapatkan dokumen yang hanya mengandung unsur marketing atau perpustakaan saja.

# LOGICAL OR

- Memperbolehkan penelusur untuk secara spesifik menggunakan alternatif *term* (atau konsep) yang mengindikasikan dua konsep sesuai dengan tujuan penelusuran. Hal ini menjadikan hasil penelusuran menjadi lebih luas, karena adanya alternatif dalam pernyataan query.
- Formulasi pernyataan sederhana : A OR B.
- Contoh :  
marketing OR library

Dengan query tersebut maka kita akan mendapatkan dokumen yang mengandung unsur marketing saja, perpustakaan saja, atau yang mengandung unsur marketing dan perpustakaan.

# LOGICAL NOT

- Dapat mengecualikan item-item dari seperangkat term penelusuran.
- Pernyataan formulasi sederhana :

A NOT B

- Contoh :

marketing NOT library

Ini artinya kita hanya menginginkan dokumen yang mengandung unsur marketing yang di dalamnya tidak ada unsur perpustakaan..

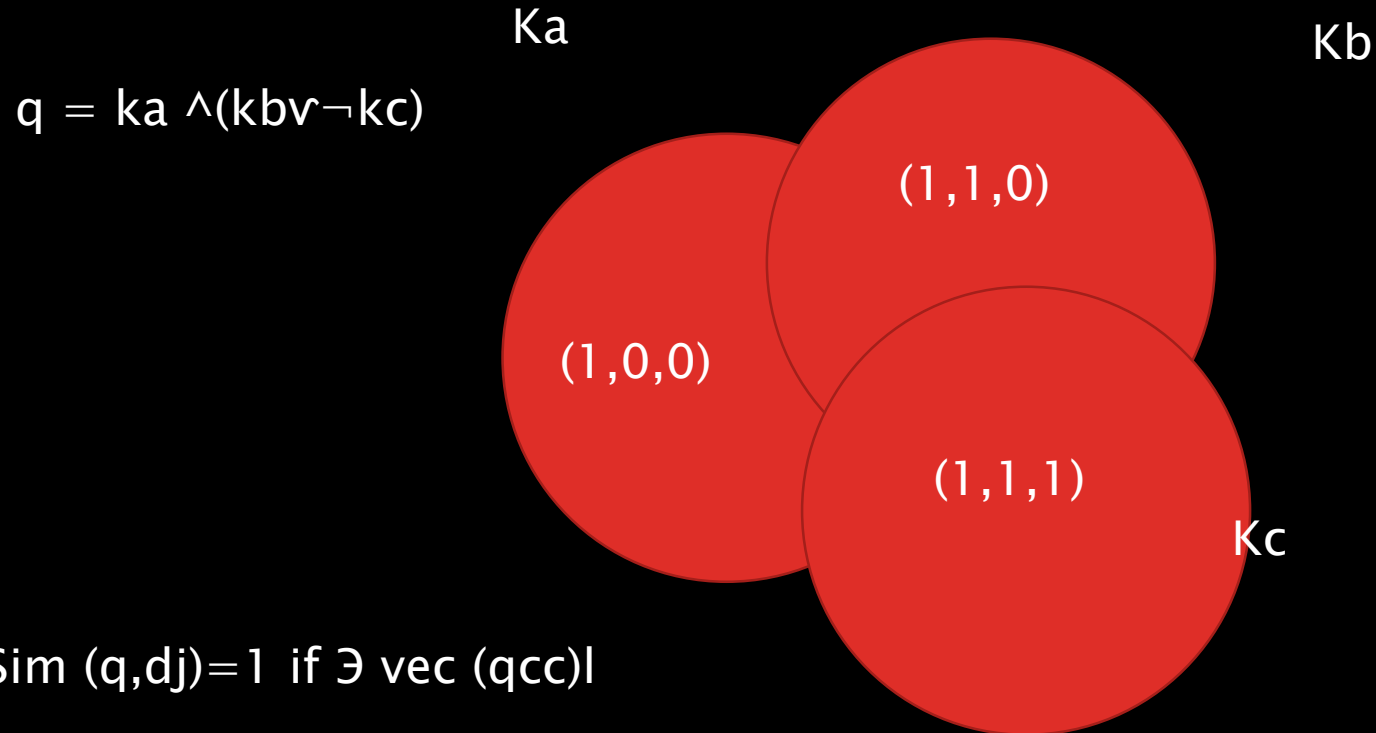
# KOMBINASI LOGICAL AND, OR, NOT

- Dapat mengkombinasikan satu pernyataan ke dalam penelusuran yang kompleks.
- Contoh :  
marketing AND library OR information centre NOT profit organization

Artinya kita ingin mendapatkan dokumen yang mengandung untur marketing dan perpustakaan tanpa unsur pusat informasi dan bukan untuk organisasi non profit.



# DIAGRAM VEN UNTUK MODEL BOOLEAN



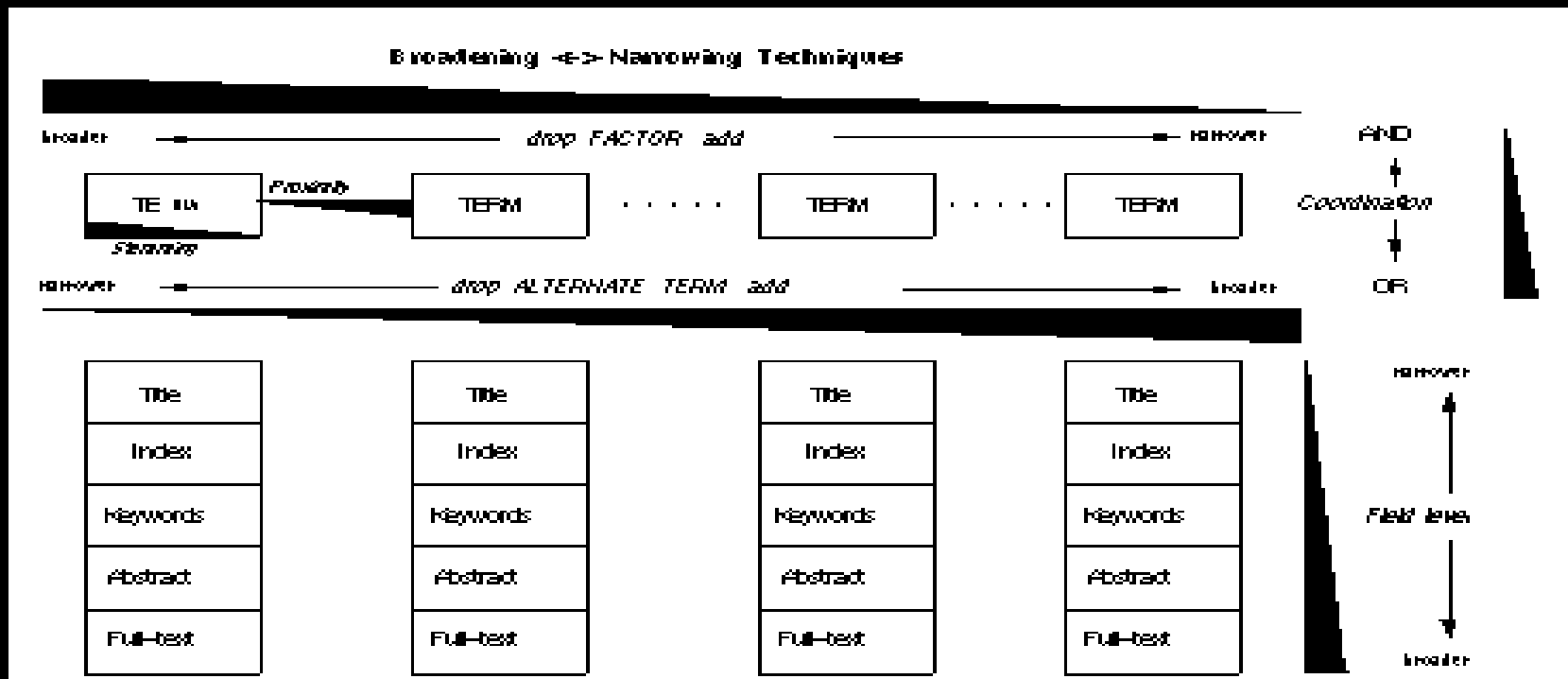
$$q = ka \wedge (kb \vee \neg kc)$$

$$\text{Sim}(q, dj) = 1 \text{ if } \exists \text{vec}(qcc) \text{ l}$$

$$\begin{aligned} &(\text{vec}(qcc) \in \text{vec}(qdnf)) \wedge \\ &(\forall ki, \\ &gi(\text{vec}(dj)) = gi(\text{vec}(qcc))) \\ &0 \text{ otherwise} \end{aligned}$$

Baeza-Yates & Ribeiro  
Neto

# TEKNIK MEMPERLUAS DAN MEMPERSEMPIT



Teknik Memperluas dan Mempersempit Hasil Penelusuran Menggunakan Boolean Logic  
<http://comminfo.rutgers.edu>

# SMART BOOLEAN

|                                    | Smart Boolean   |                         |                           |                       |                                  |                       |           |                       |              |                                    |                             |
|------------------------------------|---|-------------------------|---------------------------|-----------------------|----------------------------------|-----------------------|-----------|-----------------------|--------------|------------------------------------|-----------------------------|
| <b>Goal</b>                        | <ul style="list-style-type: none"> <li>• Structure search (re-)formulation process.</li> <li>• Use structural and contextual knowledge-bases and clarity of Boolean expressions.</li> </ul>   |                         |                           |                       |                                  |                       |           |                       |              |                                    |                             |
| <b>Methods</b>                     | <ul style="list-style-type: none"> <li>• Natural language statement is automatically translated into Boolean Topic Representation</li> <li>• Boolean Topic Representation:               <table border="0" data-bbox="517 539 1657 675"> <tr> <td>ANDs of ORs of concepts</td><td>Keyword /stem, all fields</td></tr> <tr> <td>• Conceptual info. -&gt;</td><td>Coordination and Add/Drop Factor</td></tr> <tr> <td>• Contextual info. -&gt;</td><td>Proximity</td></tr> <tr> <td>• Structural info. -&gt;</td><td>Field levels</td></tr> <tr> <td>• Synonym or word relationships -&gt;</td><td>Stemming/Truncation overlap</td></tr> </table> <p>=&gt; all this information can be used to rank documents</p> </li> <li>• Techniques to Broaden and Narrow query</li> </ul> | ANDs of ORs of concepts | Keyword /stem, all fields | • Conceptual info. -> | Coordination and Add/Drop Factor | • Contextual info. -> | Proximity | • Structural info. -> | Field levels | • Synonym or word relationships -> | Stemming/Truncation overlap |
| ANDs of ORs of concepts            | Keyword /stem, all fields   |                         |                           |                       |                                  |                       |           |                       |              |                                    |                             |
| • Conceptual info. ->              | Coordination and Add/Drop Factor  |                         |                           |                       |                                  |                       |           |                       |              |                                    |                             |
| • Contextual info. ->              | Proximity   |                         |                           |                       |                                  |                       |           |                       |              |                                    |                             |
| • Structural info. ->              | Field levels  |                         |                           |                       |                                  |                       |           |                       |              |                                    |                             |
| • Synonym or word relationships -> | Stemming/Truncation overlap   |                         |                           |                       |                                  |                       |           |                       |              |                                    |                             |
| <b>(+)</b>                         | <ul style="list-style-type: none"> <li>• No need for Boolean operators<br/>=&gt; Convert operator-free statement into ANDs of ORs</li> <li>• Assist user in query (re)formulation:<br/>by asking users targeted questions to automatically modify the query.</li> <li>• "Why irrelevant?" -&gt; activates narrowing methods.</li> <li>• "Broaden by Dropping Factors" to estimate recall.</li> </ul>  |                         |                           |                       |                                  |                       |           |                       |              |                                    |                             |
| <b>(-)</b>                         | <ul style="list-style-type: none"> <li>• How to visualize ?               <ul style="list-style-type: none"> <li>• Conceptual query representation (BTR)</li> <li>• Query modification techniques and their effects</li> <li>• Structured relevance feedback</li> </ul> </li> </ul>   |                         |                           |                       |                                  |                       |           |                       |              |                                    |                             |

<http://comminfo.rutgers.edu>

# MODEL EXTENDED BOOLEAN

|         | Extended Boolean Models  |
|---------|--|
| Goal    | <ul style="list-style-type: none"> <li>• Less strict Boolean operators</li> <li>• Ranked output</li> </ul>   |
| Methods | <div> <input checked="" type="checkbox"/> </div> <ul style="list-style-type: none"> <li>• Fuzzy logic</li> </ul> <p>[OR <math>\rightarrow</math> max], [AND <math>\rightarrow</math> min] and [NOT <math>\rightarrow</math> 1 - max]</p> <p>(-) Lack of sensitivity of min and max:<br/> <math>\min(0.2, 0.8) = \min(0.2, 0.3)</math>.</p> |

Model Perluasan Boolean digunakan untuk :

1. Apabila operator Boolean terlalu ketat dan perlu diperhalus.
2. Pendekatan standar Boolean tidak menyediakan fitur ranking sehingga pendekatan dan metode pendeskripsian dapat membantu perankingan dokumen yang relevan [Fox and Koll 1988, Marcus 1991].
3. Model Boolean tidak mendukung tugas pertimbangan query atau term dari sebuah dokumen.

<http://comminfo.rutgers.edu>

# VECTOR DAN PROBABILISTIK

| Statistical       | Vector Space   | Probabilistic                                |
|-------------------|--|--|
| <b>Motivation</b> | Simplify query formulation<br>Ability to control output  | Address uncertainty in query representations |
| <b>Goal</b>       | Rank the output based on<br>Similarity                      Probability of Relevance   |  |
| <b>Methods</b>    | Cosine measure   | Use of different models                      |
| <b>Source</b>     | <p><b>Query Term Statistics</b></p> <p><u>Vector-Space:</u></p> <ul style="list-style-type: none"> <li>• <math>\text{similarity}(Q,D) = \sum (w_{iq} \times w_{ij}) / \text{"normalizer"}</math><br/>             where <math>w_{iq} = (0.5 + 0.5 \text{ freq}_{iq} / \text{maxfreq}_q) \times \text{idf}(i)</math><br/> <math>w_{ij} = \text{freq}_{ij} \times \text{idf}(i)</math></li> <li>• inverse term freq. in collection <math>\text{idf}(i) = \log_2 (N - n(i)) / n(i)</math>.</li> </ul> <p><u>Probabilistic:</u></p> <ul style="list-style-type: none"> <li>• term weight <math>= \log [(r_i / R - r_i) / ((n_i - r_i) / ((N - n_i) - (R - r_i)))]</math><br/> <math>= \text{"(hits / misses) / (false alarms / correct misses)"}"</math></li> <li>• <math>\text{similarity}_{jk} = \sum (C + \text{idf}(i)) \times \text{tf}(i,j)</math><br/>             where <math>\text{tf}(i,j) = K + (1-K) (\text{freq}(i,j) / \text{maxfreq}(j))</math>.</li> </ul> |  |
| <b>Issues</b>     | <ul style="list-style-type: none"> <li>• How to express NOT ?</li> <li>• Proximity searches ?</li> <li>• Limited expressive power</li> <li>• Computationally intensive</li> <li>• Assumes that terms are independent.</li> <li>• Lack of structure to represent important linguistic features</li> <li>• How to better visualize the retrieved set ?</li> </ul>  |  |

<http://comminfo.rutgers.edu>

# VECTOR QUERY

- Koleksi dokumen  $n$  dengan perbedaan term  $t$  dapat direpresentasikan dengan sebuah matrix.

|          | $T_1$    | $T_2$    | .... | $T_t$    |
|----------|----------|----------|------|----------|
| $D_1$    | $w_{11}$ | $w_{21}$ | ...  | $w_{t1}$ |
| $D_2$    | $w_{12}$ | $w_{22}$ | ...  | $w_{t2}$ |
| $\vdots$ | $\vdots$ | $\vdots$ |      | $\vdots$ |
| $\vdots$ | $\vdots$ | $\vdots$ |      | $\vdots$ |
| $D_n$    | $w_{1n}$ | $w_{2n}$ | ...  | $w_{tn}$ |

- Sebuah query juga dapat direpresentasikan sebagai sebuah vector seperti sebuah dokumen.



# LATIHAN

- Temu kembali informasi dalam sebuah database dengan menggunakan operator boolean logic operators.

Contoh, kunjungi <http://online.sagepub.com> dan telusur beberapa tema tugas perkuliahan.



# Terima Kasih

Contact Me @nuningkurniasih

