# Applying the User-over-Ranking Hypothesis to Query Formulation

Matthias Hagen     Benno Stein

Bauhaus-Universität Weimar
`matthias.hagen@uni-weimar.de`

ICTIR 2011
Bertinoro, Italy
September 14, 2011

What is the User-over-Ranking hypothesis?

Queries returning as many results as the user can consider increase retrieval performance.

Queries returning as many results as the user can consider increase retrieval performance.

Fine print: If ranking works: great!
Use case is not some query like ebay.
But more involved information needs,
automatic systems, etc.

# Assumption 1: More keywords = more specific



**Specificity of Queries**

**Specificity of Queries**

**Specificity of Queries**

**Specificity of Queries**

Specificity of Queries

Probability for Retrieval Success

What is this hypothesis good for?

# Query Formulation

## Scenario

- Given a set $W$ of keywords
- Find a good query $Q \subseteq W$

# Query Formulation

## Scenario

- Given a set $W$ of keywords
- Find a good query $Q \subseteq W$

## Previous approach [Lee et al., CIKM 2009]

Learnt ranking function identifies the $m$ best keywords from $W$.

Based on:

- Known relevant documents
- Unrestricted index access
- Manually tuned $m$ for each set $W$

## Scenario

- User accessed a document
- But did not store it

## Scenario

- User accessed a document
- But did not store it

How can she find it again?

## Scenario

- User accessed a document
- But did not store it

How can she find it again?

## Solution

- Remember some keywords

`information retrieval    query formulation`
`web search    search session    user support`
`search engine    cost optimization`

- Query a search engine

But what query to formulate with the keywords?

information retrieval ~~query formulation~~

~~web search~~ ~~search session~~ ~~user support~~

~~search engine~~ ~~cost optimization~~

# Single keywords?

~~information retrieval~~    query formulation

~~web search~~    ~~search session~~    ~~user support~~

~~search engine~~    ~~user optimization~~

# Single keywords?

~~information/retrieval~~     ~~query/formulation~~

web search     ~~search/session~~     ~~user/support~~

~~search/engine~~     ~~cost/optimization~~

~~information retrieval~~ ~~query formulation~~

web search ~~search session~~ ~~user support~~

~~search engine~~ ~~cost optimization~~

information retrieval    query formulation

web search    search session    user support

search engine    cost optimization

`information retrieval`   `query formulation`

`web search`   `search session`   `user support`

`search engine`   `cost optimization`

Remember the hypothesis . . . not too many results!

information retrieval      query formulation

web search      search session      ~~user support~~

search engine      cost optimization

# "As many keywords as possible"-Query

## Characteristics

- Captures most of the remembered keywords
- Best possible description of the known-item
- Not too many results $\rightarrow$ user can check complete list

# "As many keywords as possible"-Query

## Characteristics

- Captures most of the remembered keywords
- Best possible description of the known-item
- Not too many results  $\rightarrow$  user can check complete list

## Problem

- Relevant documents not known
- No web index at user site
- Query size not known

# "As many keywords as possible"-Query

## Characteristics

- Captures most of the remembered keywords
- Best possible description of the known-item
- Not too many results $\rightarrow$ user can check complete list

## Problem

- Relevant documents not known
- No web index at user site $\rightarrow$ Lee et al. not applicable
- Query size not known

We propose an approach for this scenario ...

# Problem Statement with Capacity

## PROMISING QUERY

- Given:
  1. A set $W$ of keywords
  2. An upper bound $k$ on the result list length

- Find a largest query $Q \subseteq W$ yielding at most $k$ results

## Optimization Problem!

Minimize the number of submitted web queries to find $Q$.

# Problem Statement with Capacity

## PROMISING QUERY

- Given:
  1. A set $W$ of keywords
  2. An upper bound $k$ on the result list length

- Find a largest query $Q \subseteq W$ yielding at most $k$ results

## Optimization Problem!

Minimize the number of submitted web queries to find $Q$.

$w_1, w_2, w_3, w_4$

$w_1, w_2, w_3, w_4$

$w_1$

$w_1, w_2, w_3, w_4$

$w_1$ [10 000 results]
(> 100)

$w_1, w_2, w_3, w_4$

$w_1$ [10 000 results]
(> 100)

$w_1 \wedge w_2$

$w_1, w_2, w_3, w_4$

$w_1$ [10 000 results]
(> 100)

$w_1 \wedge w_2$ [0 results]

$w_1, w_2, w_3, w_4$

$w_1$ [10 000 results]
(> 100)

$w_1 \wedge w_2$ [0 results]

$w_1 \wedge w_3$ [90 results]
(< 100)

$w_1, w_2, w_3, w_4$

$w_1$ [10 000 results]
(> 100)

$w_2$ [10 000 results]

$w_1 \wedge w_2$ [0 results]

$w_1 \wedge w_3$ [90 results]

$w_2 \wedge w_3$ [5 000 results]

$w_1 \wedge w_3 \wedge w_4$ [0 results]

$w_2 \wedge w_3 \wedge w_4$ [60 results]

$w_1, w_2, w_3, w_4$

$w_1$ [10 000 results]
(> 100)

$w_2$ [10 000 results]

$w_1 \wedge w_2$ [0 results]

$w_1 \wedge w_3$ [90 results]

$w_2 \wedge w_3$ [5 000 results]

$w_1 \wedge w_3 \wedge w_4$ [0 results]

$w_2 \wedge w_3 \wedge w_4$ [60 results]

## Major Drawback

All intermediate queries submitted.   $\rightarrow$   Bad run time!

## Major Drawback

All intermediate queries submitted.     $\rightarrow$     Bad run time!

## Idea

Estimate the result list length before query submission.

Estimate: `"information retrieval" "query formulation" + "web search"`

# Co-occurrence based Estimations

Estimate: `"information retrieval"` `"query formulation"` + `"web search"`

Known: `"information retrieval"` `"query formulation"`     87 100 results

Estimate: `"information retrieval"` `"query formulation"` `+ "web search"`

Known: `"information retrieval"` `"query formulation"`  87 100 results

`"information retrieval"`  `+ "web search"`  16 % remain

# Co-occurrence based Estimations

Estimate: `"information retrieval"` `"query formulation" + "web search"`

Known:
| | | |
|---|---|---|
| `"information retrieval"` `"query formulation"` | | 87 100 results |
| `"information retrieval"` | `+ "web search"` | 16 % remain |
| `"query formulation"` | `+ "web search"` | 22 % remain |

# Co-occurrence based Estimations

Estimate: `"information retrieval" "query formulation" + "web search"`

Known:
| | | |
|---|---|---|
| `"information retrieval" "query formulation"` | | 87 100 results |
| `"information retrieval"` + `"web search"` | | 16 % remain |
| `"query formulation"` + `"web search"` | | 22 % remain |

Our estimation scheme:

$$\text{avg}(16\%\,,\,22\%) = 19\%$$
$$87\,100 \cdot 0.19 = 16\,500 \text{ results}$$

## Co-occurrence based Estimations

Estimate: `"information retrieval" "query formulation" + "web search"`

| Known: | `"information retrieval" "query formulation"` | 87 100 results |
| | `"information retrieval"`      `+ "web search"` | 16 % remain |
| | `"query formulation"`      `+ "web search"` | 22 % remain |

Our estimation scheme:

$$\mathrm{avg}(16\%\,,\,22\%) = 19\%$$

$$87\,100 \cdot 0.19 = 16\,500 \text{ results}$$

Control:     Google     35 700 results

## Co-occurrence based Estimations

Estimate: `"information retrieval" "query formulation" + "web search"`

Known:

| `"information retrieval" "query formulation"` | 87 100 results |
| `"information retrieval"` + `"web search"` | 16 % remain |
| `"query formulation"` + `"web search"` | 22 % remain |

Our estimation scheme:

$$\text{avg}(16\% \, , \, 22\%) \; = \; 19\%$$
$$87\,100 \cdot 0.19 \; = \; 16\,500 \text{ results}$$

Control:          Google          35 700 results

### Observation

Our scheme usually underestimates the real result list length.

$w_1, w_2, w_3, w_4$

$w_1, w_2, w_3, w_4$

$w_1$

$w_1, w_2, w_3, w_4$

$w_1$ [15 000 estimated]

$w_1, w_2, w_3, w_4$

$w_1$ [15 000 estimated]

$w_1 \wedge w_2$ [50 estimated]

$w_1, w_2, w_3, w_4$

$w_1$ [15 000 estimated]

~~$w_1 \wedge w_2$~~ [50 estimated]
[0 results]

$w_1, w_2, w_3, w_4$

$W_1$ [15 000 estimated]

~~$W_1 \wedge W_2$~~ [50 estimated]
[0 results]

$W_1 \wedge W_3$ [70 estimated]

$w_1, w_2, w_3, w_4$

$W_1$ [15 000 estimated]

$W_1 \wedge W_2$ [50 estimated] [0 results]      $W_1 \wedge W_3$ [70 estimated] [90 results]

$w_1, w_2, w_3, w_4$

$w_1$ [15 000 estimated]

$w_1 \wedge w_2$ [50 estimated]
[0 results]

$w_1 \wedge w_3$ [70 estimated]
[90 results]
(< 100)

$w_1, w_2, w_3, w_4$

$w_1$ [15 000 estimated]

$w_1 \wedge w_2$ [50 estimated]
[0 results]

$w_1 \wedge w_3$ [70 estimated]
[90 results]
(< 100)

$w_1 \wedge w_3 \wedge w_4$ [30 estimated]

$w_1, w_2, w_3, w_4$

$w_1$ [15 000 estimated]

$w_1 \wedge w_2$ [50 estimated]
[0 results]

$w_1 \wedge w_3$ [70 estimated]
[90 results]
(< 100)

$w_1 \wedge w_3 \wedge w_4$ [30 estimated]
[0 results]

# "Informed" Baseline



$w_1, w_2, w_3, w_4$

$w_1$ [15 000 estimated]

$w_1 \wedge w_2$ [50 estimated]
[0 results]

$w_1 \wedge w_3$ [70 estimated]
[90 results]
(< 100)

$w_1 \wedge w_3 \wedge w_4$ [30 estimated]
[0 results]

$w_1, w_2, w_3, w_4$

$W_1$ [15 000 estimated]

$W_2$ [20 000 estimated]

$W_1 \wedge W_2$ [50 estimated]
[0 results]

$W_1 \wedge W_3$ [70 estimated]
[90 results]
(< 100)

$W_2 \wedge W_3$ [1 500 estimated]

$W_1 \wedge W_3 \wedge W_4$ [30 estimated]
[0 results]

$W_2 \wedge W_3 \wedge W_4$ [40 estimated]

$w_1, w_2, w_3, w_4$

$W_1$ [15 000 estimated]

$W_2$ [20 000 estimated]

$W_1 \wedge W_2$ [50 estimated] [0 results]

$W_1 \wedge W_3$ [70 estimated] [90 results]

$W_2 \wedge W_3$ [1 500 estimated]

$W_1 \wedge W_3 \wedge W_4$ [30 estimated] [0 results]

$W_2 \wedge W_3 \wedge W_4$ [40 estimated] [60 results]

Informed baseline $+$ heuristic reordering of the keywords at each step

# Experimental Setup

## Corpus

- 775 papers on Computer Science (the known-items)
- 15 keywords extracted from each

## System

- Bing API as search engine
- Set $k = 100$
- Measure number of submitted Web queries

# Experimental Setup

## Corpus

- 775 papers on Computer Science (the known-items)
- 15 keywords extracted from each

## System

- Bing API as search engine
- Set $k = 100$
- Measure number of submitted Web queries

# Experimental Results

| Number of keywords | | 5 | 10 | 15 |
|---|---|---|---|---|
| Promising query | not possible | 614 | 328 | 86 |
| | found | 161 | 447 | 689 |
| Avg. queries submitted | heuristic | 10.39 | **24.93** | **53.78** |
| | informed | **10.36** | 27.01 | 108.78 |
| | baseline | 11.81 | 30.94 | 116.22 |

Almost the end: The take-away messages!

## Results

- User-over-Ranking
  - longer queries → fewer results
  - optimum retrieval performance
    → user capacity
- Heuristic for promising queries
- Use cases:
  - Known-item finding
  - Empty results lists
  - Query sessions

## Future Work

- Co-occurrence source
- User study

# What we have (not) done

## Results

- User-over-Ranking
  - longer queries $\rightarrow$ fewer results
  - optimum retrieval performance
    $\rightarrow$ user capacity
- Heuristic for promising queries
- Use cases:
  - Known-item finding
  - Empty results lists
  - Query sessions

## Future Work

- Co-occurrence source
- User study

## Results

- User-over-Ranking
  - longer queries $\rightarrow$ fewer results
  - optimum retrieval performance
    $\rightarrow$ user capacity
- Heuristic for promising queries
- Use cases:
  - Known-item finding
  - Empty results lists
  - Query sessions

## Future Work

- Co-occurrence source
- User study

# Thank you
☺