# Introducing the User-over-Ranking hypothesis

Benno Stein    Matthias Hagen

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

ECIR 2011
Dublin, Ireland
April 21, 2011

# User-over-Ranking stated informally

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

# User-over-Ranking stated informally

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

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

# Assumption 1: More keywords = more specific

# Assumption 1: More keywords = more specific

# Assumption 1: More keywords = more specific



**Specificity of Queries**

# Assumption 2: User can arbitrarily specify information need



**Specificity of Queries**

# Assumption 2: User can arbitrarily specify information need



**Specificity of Queries**

# Assumption 3: User can consider about $k$ results.



Specificity of Queries

# Hypothesis: Specificity matches $k$ = Optimum retrieval



**Specificity of Queries**

**Probability for Retrieval Success**

What about empirical evidence?

# Experimental Setting: AOL log

- Cleaning (bots, URL queries, encoding problems)
- Query duplicates removed
- 4.4 million unique queries ($\leq 22$ keywords)
- Submitted to Bing API
- Result list length estimates stored

# AOL log result list length distribution in 3D

# Median AOL log result list length in 2D

# Experimental Setting: TREC Robust04

- 530 000 newswire documents
- BM25 indexed with Terrier
- Nounphrase extraction for TREC topics 301–450, 601–700
- Submitted all combinations to Terrier
- Result lists stored

- Assumed capacity $k = 100$

# Avg. NDCG@100 per result list length (Robust04)

Almost the end: The take-away messages!

# What we have done

## Results

- When ranking works: fine!
- Else: User-over-Ranking
  - longer queries $\rightarrow$ fewer results
  - optimum retrieval performance
    $\rightarrow$ user capacity
- Empirical evidence

## Future Work

- Apply hypothesis to query formulation

# What we have (not) done

### Results

- When ranking works: fine!
- Else: User-over-Ranking
  - longer queries $\rightarrow$ fewer results
  - optimum retrieval performance
    $\rightarrow$ user capacity
- Empirical evidence

### Future Work

- Apply hypothesis to query formulation

# What we have (not) done

## Results

- When ranking works: fine!
- Else: User-over-Ranking
  - longer queries $\rightarrow$ fewer results
  - optimum retrieval performance
    - $\rightarrow$ user capacity
- Empirical evidence

## Future Work

- Apply hypothesis to query formulation

# Thank you

☺