Meta-Evaluating Search Engines

Frameworks for Evaluating and Meta-Evaluating Search Engines

Ph.D Thesis Defense

Peter B. Golbus Advisor: Javed A. Aslam

Northeastern University College of Computer and Information Science

June 19th. 2014

Golbus (NEU CCIS)

Meta-Evaluating Search Engines

-June 19th, 2014 1 / 82

- ∢ ⊢⊒ →

Evaluation: The Bottleneck in IR Research

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework RIC Conditional Rank Correlatic

IE

Classifier Metasearch

Conclusion

- Web search, Diversification, Topic Distillation, Knowledge Acquisition, Temporal Summarization, Legal search, Enterprise search, Microblog search...
- Little to no formal problem description

Evaluation: The Bottleneck in IR Research

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework RIC Conditional Rank Correlatio

ID

Classifier Metasearch

Conclusion

- Web search, Diversification, Topic Distillation, Knowledge Acquisition, Temporal Summarization, Legal search, Enterprise search, Microblog search...
- Little to no formal problem description
- We don't really know what we're trying to do, so...

Evaluation: The Bottleneck in IR Research

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

- Framework RIC Conditional Rank Correlatic
- ID
- Classifier Metasearch
- Conclusion

- Web search, Diversification, Topic Distillation, Knowledge Acquisition, Temporal Summarization, Legal search, Enterprise search, Microblog search...
- Little to no formal problem description
- We don't really know what we're trying to do, so...
 - 1 Are we improving the right thing?
 - 2 How do we know when we've done it?

Outline

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework RIC Conditional Rank Correlatio

ID

Classifier Metasearch

Conclusion

- **1** Motivation: Evaluating Information Retrieval
- 2 Targeted Meta-Evaluation for Understanding Diversity

3 Probabilistic Framework for Evaluation and Rank Correlation

4 Information Difference



Outline

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework RIC Conditional Rank Correlatio

ID

Classifier Metasearcl

Conclusion

1 Motivation: Evaluating Information Retrieval

Targeted Meta-Evaluation for Understanding Diversity

Probabilistic Framework for Evaluation and Rank Correlation

Information Difference



< E

Image: A math a math

The Ad Hoc Retrieval Task

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework RIC Conditional Rank Correlation

ID

Classifier Metasearch

Conclusion

- Information need expressed in detailed narrative *e.g.* Documents discussing electrical output of the Three Gorges Dam
 - Known as topics or queries
- Find documents containing even one pertinent fact
 - Set of "relevance labels" known as a QREL
- Recall-oriented: user interacts with *all* ranked documents

The Ad Hoc Retrieval Task

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework RIC Conditional Rank Correlation

ID

Classifier Metasearch

Conclusion

- Information need expressed in detailed narrative *e.g.* Documents discussing electrical output of the Three Gorges Dam
 - Known as topics or queries
- Find documents containing even one pertinent fact
 - Set of "relevance labels" known as a QREL
- Recall-oriented: user interacts with *all* ranked documents

e.g. Annual NIST-sponsored Text REtrieval Conference (TREC)

Ad Hoc Evaluation



Golbus

Motivation

Diversity

```
Framework
RIC
Conditional
Rank Correlatio
```

ID

Classifier Metasearch

Conclusion



Precision: percentage of retrieved documents that are relevant

 $Prec@k = \frac{|\text{Relevant}| \cap |\text{Retrieved}|}{|\text{Retrieved}|}$

Recall: percentage of relevant documents that are retrieved

$$Rec@k = \frac{|\text{Relevant}| \cap |\text{Retreived}|}{|\text{Relevant}|}$$

Meta-Evaluating Search Engines

Ad Hoc Evaluation

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework RIC Conditional Rank Correlatio

ID

Classifier Metasearch

Conclusion

- Precision: percentage of retrieved documents that are relevant
- Recall: percentage of relevant documents that are retrieved
- Average Precision: average precision at ranks of relevant documents



Web Search

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework RIC Conditional Rank Correlatio

ID

Classifier Metasearch

Conclusion

- Information need expressed in 2 or 3 word query
- Find high quality, highly relevant documents
- Precision-oriented: user interacts with few documents

Image: Image:

Web Search

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework RIC Conditional Rank Correlatic

ID

Classifier Metasearch

Conclusion

- Information need expressed in 2 or 3 word query
- Find high quality, highly relevant documents
- Precision-oriented: user interacts with few documents
 - Usually

э

-

Image: A math a math

Web Evaluation

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework RIC Conditional Rank Correlati

IE

Classifier Metasearch

Conclusion

Find high quality, highly relevant documents

- Use graded relevance labels
- Precision-oriented: user interacts with few documents
 - Rank-based discounting
- $eval = \langle gain, discount \rangle$

Rank	Grade	Gain	Discount
1	4	15	1
2	0	0	0.5
3	2	3	0.38
4	0	0	0.33
5	3	7	0.30

$$(\mathsf{n})\mathsf{DCG} = rac{2^{grade}-1}{\log(rank+1)}$$

nDCG@5 = 18.24

Golbus (NEU CCIS)

Ad Hoc Evaluation

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework RIC Conditional Rank Correlatio

ID

Classifier Metasearch

Conclusion

- Precision: percentage of retrieved documents that are relevant
- Recall: percentage of relevant documents that are retrieved
- Average Precision: average precision at ranks of relevant documents



Web Evaluation

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework RIC Conditional Rank Correlati

IC

Classifier Metasearch

Conclusion

Find high quality, highly relevant documents

- Use graded relevance labels
- Precision-oriented: user interacts with few documents
 - Rank-based discounting
- $eval = \langle gain, discount \rangle$

Rank	Grade	Gain	Discount
1	4	15	1
2	0	0	0.5
3	2	3	0.38
4	0	0	0.33
5	3	7	0.30

$$(\mathsf{n})\mathsf{DCG} = rac{2^{grade}-1}{\log(rank+1)}$$

< 🗇 🕨

nDCG@5 = 18.24

Golbus (NEU CCIS)

э

Web Evaluation

Meta-Evaluating Search Engines

Motivation

• $eval = \langle gain, discount \rangle$

Rank	Grade	Gain	Discount	
1	4	15	1	
2	0	0	0.5	
3	2	3	0.38	(n)DCG =
4	0	0	0.33	
5	3	7	0.30	

Conclusion

nDCG@5 = 18.24

Doesn't Match Actual Behavior!*

*Golbus et al. WWW '14

- 一司

Goals

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework RIC Conditional Rank Correlatio

ID

Classifier Metasearch

Conclusion

- 1 Use targeted meta-evaluations to better understand IR problems
- 2 New, powerful, highly interpretable information-theoretic toolkit
 - Evaluation and meta-evaluation within a single unified framework

э

3

Image: A matrix of the second seco

Contributions

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework RIC Conditional Rank Correlatio

ID

Classifier Metasearch

Conclusion

- Targeted Meta-Evaluation for Understanding Diversity
 Document Selection Sensitivity
- 2 Probabilistic Framework for Evaluation and Rank Correlation
 - Information-Theoretic Evaluation Measure
 - Conditional Rank Correlation
- **3** Information Difference
 - Similarity Classifier
 - Selecting Systems for Metasearch

Outline

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework RIC Conditional Rank Correlatio

ID

Classifier Metasearc

Conclusion

Motivation: Evaluating Information Retrieval

2 Targeted Meta-Evaluation for Understanding Diversity

Probabilistic Framework for Evaluation and Rank Correlation

Information Difference



∃ → < ∃</p>

Image: A matrix and a matrix

Novelty & Diversity

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework RIC Conditional Rank Correlation

ID

Classifier Metasearch

Conclusion

Subtopics

What is the information need of the query "freddie mercury" biography, memorabilia, music?

Goal

Cover as many subtopics as possible by rank k

Golbus, Aslam & Clarke. Inf. Retr. 16(4) = > < =>

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework RIC Conditional Rank Correlatio

ID

Classifier Metasearch

Conclusion

Two Factors

1 Diversity: how well did you order the documents?

3

∃ → (∃ →

A B A B A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework RIC Conditional Rank Correlatio

ID

Classifier Metasearch

Conclusion

Two Factors

Diversity: how well did you order the documents?
 Performance: how many good documents did you find?

э

3

Image: A matrix

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework RIC Conditional Rank Correlation

ID

Classifier Metasearch

Conclusion

- Diversity measures are traditional measures with an added diversity component
 - *e.g.* Intent Aware measures: evaluate with respect to each subtopic and compute the average
 - $\blacksquare \ nDCG \rightarrow nDCG-IA$
- Are we really measuring diversity?

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework RIC Conditional Rank Correlation

ID

Classifier Metasearch

Conclusion

- Diversity measures are traditional measures with an added diversity component
 - *e.g.* Intent Aware measures: evaluate with respect to each subtopic and compute the average
 - $\blacksquare \ nDCG \rightarrow nDCG-IA$
- Are we really measuring diversity?

No!

Measure Sensitivity

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework RIC Conditional Rank Correlation

IE

Classifier Metasearch

Conclusion

How **sensitive** to changes in system performance?

 Discriminative power: percentage of system pairs that are statistically significantly different

Image: Image:

Measure Sensitivity

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework RIC Conditional Rank Correlatio

ID

Classifier Metasearch

Conclusion

How sensitive to changes in system performance?

- Discriminative power: percentage of system pairs that are statistically significantly different
- Sensitivity is necessary but not sufficient

Measure Sensitivity

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework RIC Conditional Rank Correlatio

ID

Classifier Metasearch

Conclusion

How **sensitive** to changes in system performance?

- Discriminative power: percentage of system pairs that are statistically significantly different
- Sensitivity is necessary but not sufficient

What is the measure sensitive to?

Does Document Order Matter?



Conclusion

ERR-IA	2010	2011
Actual	0.571	0.544

D#-nDCG	2010	2011
Actual	0.583	0.600

Table: Discriminative power on actual runs

Image: A matrix

э

Does Document Order Matter?

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework RIC Conditional Rank Correlatio

ID

Classifier Metasearch

Conclusion

ERR-IA	2010	2011
Actual	0.571	0.544
Permutations	0.039	0.031

D#-nDCG	2010	2011
Actual	0.583	0.600
Permutations	0.036	0.019

Table: Discriminative power on actual runs and permutations of relevant documents

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework RIC Conditional Rank Correlatio

ID

Classifier Metasearch

Conclusion

Two Factors

Diversity: how well did you order the documents?
 Performance: how many good documents did you find?

Example

biography

- memorabilia
- music
- Not relevant

э

3

Image: A matrix

Two Factors

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework RIC Conditional Rank Correlatio

ID

Classifier Metasearch

Conclusion

Diversity: how well did you order the documents?
 Performance: how many good documents did you find?

Bad performance \Rightarrow bad diversity



Image: Image:

Meta-Evaluating Search Engines

э

Two Factors

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework RIC Conditional Rank Correlatio

IE

Classifier Metasearch

Conclusion

Diversity: how well did you order the documents?
 Performance: how many good documents did you find?

Good diversity \Rightarrow good performance

Image: Image:

Golbus (NEU CCIS)

Meta-Evaluating Search Engines

June 19th, 2014 20 / 82

э

Two Factors

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework RIC Conditional Rank Correlatio

ID

Classifier Metasearch

Conclusion

Diversity: how well did you order the documents?
 Performance: how many good documents did you find?

Good performance \neq good diversity

Image: Image:

э

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework RIC Conditional Rank Correlatio

ID

Classifier Metasearch

Conclusion

Two Factors

Diversity: how well did you order the documents?
 This does not matter

Performance: how many good documents did you find?

This does

э

< E

.⊒ . ►

Image: A matrix and a matrix

Research Bottleneck Example

Meta-Evaluating Search Engines

Diversity

The Problem

We don't really know what we're trying to do, so...

- **1** Are we improving the right thing?
- 2 How do we know when we've done it?

Recent Work

- Hypothesis: Metasearch improves diversity
 - Metasearch improves performance
 - Improving performance improves diversity evaluation scores

Research Bottleneck Example

Meta-Evaluating Search Engines

The Problem

We don't really know what we're trying to do, so...

- **1** Are we improving the right thing?
- 2 How do we know when we've done it?

Recent Work

- IE
- Classifier Metasearch

Diversity

Conclusion

- Hypothesis: Metasearch improves diversity
 - Metasearch improves performance
 - Improving performance improves diversity evaluation scores

Hypothesis not tested by standard diversity evaluation paradigm!

Golbus (NEU CCIS)

Meta-Evaluating Search Engines

June 19th, 2014 22 / 82

3 🕨 🖌 3

A 🖓 h
Document Selection Sensitivity

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework RIC Conditional Rank Correlatio

ID

Classifier Metasearch

Conclusion

- How do we control for performance?
 - Use only relevant documents
- What do we want to know?
 - Does document order impact evaluation score?

Document Selection Sensitivity (dss)

- For a measure *M*, evaluate random permutations of relevant documents
- Document Selection Sensitivity is the *coefficient of* variation

$$dss(M) = rac{standard\ deviation}{mean}$$

Document Selection Sensitivity



Conclusion

- Scale-invariant measure of dispersion
- A DSS of .1 means 68% of scores are within +/- 10% of the average score (assuming normally distributed)

Improved Measures



Golbus

Motivation

Diversity

Framework RIC Conditional Rank Correlatio

ID

Classifier Metasearch

Conclusion

DSS	ERR	DCG	RBP
Cascade	0.204	0.152	0.170
D#	0.116	0.096	0.111
α #-IA	0.871	0.801	0.801

Intuition: Focus on difficult topics

• *i.e.* leverage as much diversity information from the collection as you can

Image: A matrix of the second seco

Contribution



Golbus

Motivation

Diversity

Framework RIC Conditional Rank Correlatio

IE

Classifier Metasearch

Conclusion

Targeted Meta-Evaluation for Understanding Diversity Document Selection Sensitivity

э

-

Image: A matrix and a matrix

Outline

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework

RIC Conditional Rank Correlation

ID

Classifier Metasearch

Conclusion

Motivation: Evaluating Information Retrieval

Targeted Meta-Evaluation for Understanding Diversity

Probabilistic Framework for Evaluation and Rank Correlation
 Information-Theoretic Evaluation Measure

Conditional Rank Correlation

4 Information Difference

5 Conclusion

-

Observation 1: From Geometric to Probabilistic

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework

RIC Conditional Rank Correlation

IE

Classifier Metasearch

Conclusion

Currently: Vector space model of evaluation
eval(S) = (gain, discount)

Desired: Probabilistic model

 $eval(S) = I(QREL; RANKED_LIST)$

Golbus & Aslam. SIGIR '13.

Golbus (NEU CCIS)

Meta-Evaluating Search Engines

Observation 2: Comparing Orderings



Observation 2: Comparing Orderings

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework

RIC Conditional Rank Correlation

ID

Classifier Metasearch

Conclusion

- \blacksquare Kendall's τ depends on preferences, not ranks
- QRELs and ranked lists encode document preferences
- Lets measure "rank" correlation between these preferences

Encoding Preferences

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework

RIC Conditional Rank Correlation

ID

Classifier Metasearch

Conclusion

{Highly Relevant Docs} {Relevant Docs} {Non-Relevant Docs} {Unjudged Docs}

Doc 1 Doc 2 E Doc *k* {Unretrieved Docs}

QREL

Ranked List

∃ ► < ∃</p>

< 4 ₽ × <

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework

RIC Conditional Rank Correlation

ID

Classifier Metasearcl

Conclusion

1 Sample space

2 Distribution over sample space

3 Random variables

Image: A matrix of the second seco

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework

RIC Conditional Rank Correlation

ID

Classifier Metasearch

Conclusion

Sample space

- $\Omega = \text{all } 2 \cdot \binom{n}{2}$ document pairs
- 2 Distribution over sample space
- 3 Random variables

-

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework

RIC Conditional Rank Correlation

ID

Classifier Metasearch

Conclusion

Sample space

- $\Omega = \text{all } 2 \cdot \binom{n}{2}$ document pairs
- 2 Distribution over sample space
 - $\bullet P = U$
- 3 Random variables

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework

RIC Conditional Rank Correlatior

IE

Classifier Metasearch

Conclusion

1 Sample space

- $\Omega = \text{all } 2 \cdot \binom{n}{2}$ document pairs
- 2 Distribution over sample space
 - $\bullet P = U$
- 3 Random variables
 - For a ranked list R, $X_R \colon \Omega \to \{-1, +1\}$ $X_R[(d_i, d_j)] = \begin{cases} 1 & \text{if } d_i \text{ appears before } d_j \text{ in } R. \\ -1 & \text{otherwise.} \end{cases}$ • $E[X_R] = 0$

Meta-Evaluating Search Engines

1 Sample space

Framework

• $\Omega = \text{all } 2 \cdot \binom{n}{2}$ document pairs

Distribution over sample space 2

P = U

- 3 Random variables
 - For a ranked list $R, X_R \colon \Omega \to \{-1, +1\}$ $X_R[(d_i, d_j)] = \begin{cases} 1 & \text{if } d_i \text{ appears before } d_j \text{ in } R. \\ -1 & \text{otherwise.} \end{cases}$ $E[X_R] = 0$

For two lists R and S

•
$$E[X_R \cdot X_S] = \frac{2c-2d}{2(c+d)} = \tau(R,S)$$

Golbus (NEU CCIS)

Power of Framework



Golbus

Motivation

Diversity

Framework

RIC Conditional Rank Correlatio

ID

Classifier Metasearch

Conclusion

Power of Framework

- 1 Generalizes to partial orderings
- 2 Flexibility of random variables
- 3 Information-theoretic interpretation

Image: Image:

Information Theory



Golbus (NEU CCIS)

Meta-Evaluating Search Engines

June 19th, 2014 34 / 82

- < ∃ →

A B A B A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A

æ

RIC



I(System; QREL)

Image: A matrix

Meta-Evaluating Search Engines

3

Conditional Rank Correlation



Meta-Evaluating Search Engines

June 19th, 2014

Conditional Rank Correlation



$$I(List_1; List_2 | List_3)$$

Information Difference



Golbus (NEU CCIS)

Meta-Evaluating Search Engines

June 19th, 2014 38 / 82

3

(日) (同) (三) (三)

Information Difference



3

(日) (同) (三) (三)

Outline

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework

RIC

Conditional Rank Correlation

ID

Classifier Metasearch

Conclusion

Targeted Meta-Evaluation for Understanding Diversity

3 Probabilistic Framework for Evaluation and Rank Correlation

- Information-Theoretic Evaluation Measure
- Conditional Rank Correlation

Information Difference

5 Conclusion

-

Image: A matrix and a matrix

-

RIC



Meta-Evaluating Search Engines

Evaluation Measure



Golbus

Motivation

Diversity

Framework

RIC

Conditional Rank Correlation

IE

Classifier Metasearch

Conclusion

Relevance Information Correlation

 $RIC(S) = I(R_S; Q)$

3

∃ → (∃ →

A B A B A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A

Evaluation Measure



Golbus

Motivation

Diversity

Framework

RIC

Conditional Rank Correlation

ID

Classifier Metasearcl

Conclusion

Relevance Information Correlation

 $RIC(S) = I(R_S; Q)$

 $RIC(S_1, \ldots, S_n)$ is well-defined

3

(日) (同) (三) (三)

Evaluation Measure



Golbus

Motivation

Diversity

Framework

RIC

Conditional Rank Correlation

ID

Classifier Metasearch

Conclusion

Relevance Information Correlation

 $RIC(S) = I(R_S; Q)$

 $RIC(S_1, \ldots, S_n)$ is well-defined

Upper bound on metasearch

э

< E

-

Image: A matrix and a matrix

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework

RIC

Conditional Rank Correlation

ID

Classifier Metasearch

Conclusion

Sample space

2 Distribution over sample space

- RIC@k*
- 3 Random variables

*Golbus & Aslam. SIGIR '14.

Golbus (NEU CCIS)

"Rank" of a Relevance Grade



Consider ranked lists consistent with the QREL: Highly Relevant

> Highly Relevant Relevant

Relevant Non-Relevant

Non-Relevant

- $R_g :=$ number of documents with relevance grade g
- *k_{min}* := minimum rank of a document with grade *g*
- $k_{max} := maximum rank of a document with grade g$

44 / 82

Probability of a Document

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework

RIC

Conditional Rank Correlation

IE

Classifier Metasearch

Conclusion

A rank discount function can be interpreted as a probability distribution.*

•
$$P_{DCG}(k) = \frac{1}{\log_2(k+1)} - \frac{1}{\log_2(k+2)}$$

Probability of a document



Probability of a document pair
 P(d_i, d_i) = βP(d_i)P(d_i)

*Carterette. SIGIR '11.

Golbus (NEU CCIS)

RIC vs nDCG @ 20

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework

RIC

Conditional Rank Correlation

IE

Classifier Metasearch

Conclusion



Golbus (NEU CCIS)

Meta-Evaluating Search Engines

June 19th, 2014 46 / 82

Outline

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework RIC Conditional Rank Correlation

ID

Classifier Metasearch

Conclusion

Targeted Meta-Evaluation for Understanding Diversity

3 Probabilistic Framework for Evaluation and Rank Correlation

- Information-Theoretic Evaluation Measure
- Conditional Rank Correlation

Information Difference

5 Conclusion

-

3 ×

Image: A matrix and a matrix

Conditional Rank Correlation



Meta-Evaluating Search Engines

Diversity Dominated by Performance

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework RIC Conditional Rank Correlation

ID

Classifier Metasearc

Conclusion

- Diversity measures are highly rank-correlated
- Diversity evaluation is dominated by performance
- Hypothesis: Correlation due to performance, not diversity

Diversity Dominated by Performance

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework RIC Conditional Rank Correlation

ID

Classifier Metasearc

Conclusion

- Diversity measures are highly rank-correlated
- Diversity evaluation is dominated by performance
- Hypothesis: Correlation due to performance, not diversity

How do we test this?

Rank Correlation



$$I(X_R; X_S) = \frac{1+\tau}{2} \log(1+\tau) + \frac{1-\tau}{2} \log(1-\tau)$$

Golbus (NEU CCIS)

э. June 19th, 2014

э

50 / 82

Conditional Rank Correlation



Golbus (NEU CCIS)

Meta-Evaluating Search Engines

June 19th, 2014 51 / 82
Comparing Evaluation Measures



Conclusion

	Measure 1	Measure 2	
1.	Bing	Google	
2.	Yandex	Bing	
3.	Google	Yandex	
4.	Baidu	Baidu	

(日) (同) (三) (三)

э

Comparing Evaluation Measures



	Measure 1	Measure 2	
1.	Bing	Google	
2.	Yandex	Bing	
3.	Google	Yandex	
4.	Baidu	Baidu	

Concordant vs. Discordant

3

-

A B A B A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A





- Performance and Diversity
- Just Diversity

Golbus (NEU CCIS)

June 19th, 2014 55 / 82

3

э.

A B A B A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A

-

-

Meta- Evaluating Search Engines			
Golbus			
Motivation			
Diversity			
Framework			
RIC Conditional			

ID

Classifier Metasearch

Conclusion

	TREC 2010	TREC 2011
$ au_I(ERR ext{-IA} ; D\# ext{-nDCG})$	0.64	0.55
$ au_I(ERR\text{-}IA \ ; D\#\text{-}nDCG \ \ nDCG)$	0.30	0.17
$ au_I(ERR ext{-}IA \ ; D\# ext{-}nDCG \ \ ERR)$	0.12	0.14
$ au_I(ERR\text{-}IA\ ;\ D\#\text{-}nDCG\ \ nDCG,\ ERR)$	0.12	0.10

3

<ロ> (日) (日) (日) (日) (日)



э

(日) (同) (三) (三)



Rank-correlation in diversity measures not due to diversity!

Golbus (NEU CCIS)

Meta-Evaluating Search Engines

June 19th, 2014 57 / 82

Image: Image:

Contribution

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework RIC Conditional Rank Correlation

ID

Classifier Metasearch

Conclusion

Probabilistic Framework for Evaluation and Rank Correlation

- Evaluation and meta-evaluation
- Information-theoretic interpretation
- Many novel applications

Image: A matrix

э

Outline

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework RIC Conditional Rank Correlatio

ID

Classifier Metasearch

Conclusion

1 Motivation: Evaluating Information Retrieval

2 Targeted Meta-Evaluation for Understanding Diversity

Probabilistic Framework for Evaluation and Rank Correlation

4 Information Difference

- Similarity Classifier
- Selecting Systems for Metasearch

5 Conclusion

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework RIC Conditional Rank Correlation

ID

Classifier Metasearch

Conclusion

How do you know when two systems are different? Measure performance delta

4

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework RIC Conditional Rank Correlatic

ID

Classifier Metasearch

Conclusion

How do you know when two systems are different?

Measure performance delta

What if their performance is the same?

1 They are essentially the same

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework RIC Conditional Rank Correlatio

ID

Classifier Metasearch

Conclusion

How do you know when two systems are different?

Measure performance delta

What if their performance is the same?

- 1 They are essentially the same
- 2 Documents with the same grades appear at the same ranks

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework RIC Conditional Rank Correlatio

ID

Classifier Metasearch

Conclusion

How do you know when two systems are different?

Measure performance delta

What if their performance is the same?

1 They are essentially the same

2 Documents with the same grades appear at the same ranks

Want to compare system behavior not system performance

Information Difference



Golbus (NEU CCIS)

Meta-Evaluating Search Engines

June 19th, 2014 61 / 82

3

(日) (同) (三) (三)

Outline

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework RIC Conditional Rank Correlation

ID

Classifier Metasearch

Conclusion

Targeted Meta-Evaluation for Understanding Diversity

3 Probabilistic Framework for Evaluation and Rank Correlation

4 Information Difference

- Similarity Classifier
- Selecting Systems for Metasearch

5 Conclusion

"Similarity"

Meta-Evaluating Search Engines

Golbus

- Motivation
- Diversity
- Framework RIC Conditional Rank Correlatio
- ID
- Classifier Metasearch
- Conclusion

- Research groups submit multiple systems
- Sort systems into bins by performance
- In each bin, mark systems from same group "similar"

Image: A matrix

ID vs Performance Delta



э

э.

A B A B A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A

Other Baselines

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework RIC Conditional Rank Correlatio

ID

Classifier Metasearch

Conclusion

1 Mutual Information: $I(R_{S_1}; R_{S_2})$

- Order information, no relevance judgments
- Base case: Kendall's τ
- **2** Jaccard Coefficient: $\frac{S_1 \cap S_2}{S_1 \cup S_2}$
 - No order information, no relevance judgments

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework RIC Conditional Rank Correlatio

ID

Classifier Metasearch

Conclusion



3



June 19th, 2014 67 / 82

A B A B A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A

э

Outline

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework RIC Conditional Rank Correlati

ID

Classifier Metasearch

Conclusion

Targeted Meta-Evaluation for Understanding Diversity

3 Probabilistic Framework for Evaluation and Rank Correlation

4 Information Difference

- Similarity Classifier
- Selecting Systems for Metasearch

5 Conclusion

• 4 ∰ ▶ 4 ∃

Metasearch

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework RIC Conditional Rank Correlatic

ID

Classifier Metasearch

Conclusion

Metasearch: combine multiple search engines into a single ranked list that is better than the worst input system.

-

Image: Image:

э

Which Systems?



Conclusion

ም.

Which Systems?



Related work: tail not max How to choose best without relevance judgments?

Golbus (NEU CCIS)

Meta-Evaluating Search Engines

June 19th, 2014 70 / 82

Which Systems?



Related work: tail not max How to choose best without relevance judgments? Information difference framework

Golbus (NEU CCIS)

June 19th, 2014 70 / 82

Approach

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework RIC Conditional Rank Correlatio

ID

Metasearch

Conclusion

Common wisdom: best, most dissimilar systems

1 How to define best and most dissimilar?

- Best: "gold standard" rankers
- Similarity: information difference framework

2 How to merge best and most dissimilar?

Facilities Location Analysis



Motivatior

Diversity

Framework RIC Conditional Rank Correlatio

ID

Metasearch

Conclusion



Synthetic data with 4 clusters

Golbus (NEU CCIS)

Meta-Evaluating Search Engines

June 19th, 2014 72 / 82

Facilities Location Analysis

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework RIC Conditional Rank Correlatio

ID

Classifier Metasearch

Conclusion



Obnoxious Facilities Dispersion

Maximal Marginal Relevance—Most unique*

Desirable Facilities Placement

<u>k-Medoids Clustering</u>—Most representative

*Carbonell & Goldstein. SIGIR '98.

Meta-Evaluating Search Engines

June 19th, 2014 73 / 82

Facilities Location Analysis

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework RIC Conditional Rank Correlatio

ID

Classifier Metasearch

Conclusion



Obnoxious Facilities Dispersion

Maximal Marginal Relevance—Most unique

Desirable Facilities Placement

<u>k-Medoids Clustering</u>—Most representative*

*Zuccon et al. ECIR '12

Golbus (NEU CCIS)

Meta-Evaluating Search Engines

June 19th, 2014 73 / 82

Selecting Systems Without Relevance Judgments





Best, most representative systems

Golbus (NEU CCIS)

Contribution

 Meta-Evaluating Search Engines
 Information Difference

 Motivation
 Information Difference

 RIC Conditional Rank Correlation
 Meta-evaluation tool for comparing systems by behavior.

Metasearch

Conclusion

3

< E

A B A B A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A

Outline

Meta-Evaluating Search Engines

Golbus

Motivatior

Diversity

Framework RIC Conditional Rank Correlatio

ID

Classifier Metasearch

Conclusion

Motivation: Evaluating Information Retrieval

Targeted Meta-Evaluation for Understanding Diversity

Probabilistic Framework for Evaluation and Rank Correlation

Information Difference



∃ → < ∃</p>

Image: A matrix and a matrix

Contributions

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework RIC Conditional Rank Correlatio

ID

Classifier Metasearch

Conclusion

- Targeted Meta-Evaluation for Understanding Diversity
 Document Selection Sensitivity
- 2 Probabilistic Framework for Evaluation and Rank Correlation
 - Information-Theoretic Evaluation Measure
 - Conditional Rank Correlation
- **3** Information Difference
 - Similarity Classifier
 - Selecting Systems for Metasearch

ID for measures

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework RIC Conditional Rank Correlatio

IE

Classifier Metasearch

Conclusion

- Information difference: IR specific rank correlation between ranked lists of *documents*
- Can we do this for ranked lists of systems, i.e. evaluation measures?

э

Image: A matrix

ID for measures

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework RIC Conditional Rank Correlatio

ID

Classifier Metasearch

Conclusion

- Information difference: IR specific rank correlation between ranked lists of *documents*
- Can we do this for ranked lists of systems, i.e. evaluation measures?
- What is the analogue of a QREL?

ID for measures

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework RIC Conditional Rank Correlatio

ID

Classifier Metasearch

Conclusion

- Information difference: IR specific rank correlation between ranked lists of *documents*
- Can we do this for ranked lists of systems, i.e. evaluation measures?
- What is the analogue of a QREL?
- Potential application: Comparing crowdsourced QRELs
Mixed Relevance Information

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework RIC Conditional Rank Correlatio

ID

Classifier Metasearch

Conclusion

Relevance Judgments

- Experts vs crowdworkers
- Grades vs preferences
 - Direct vs observed
- Estimates from rankers

Mixing Relevance Judgments

Given n QRELs, how do we use all of them?

• Currently, use evaluation pipeline: e.g. experts then clicks Our framework: $I(S; Q_1, ..., Q_n)$

- Evaluate with respect to all QRELs simultaneously
- Each QREL weighted appropriately

	Thank You
Meta- Evaluating Search Engines Golbus	
Motivation Diversity	
Framework RIC Conditional Rank Correlation	Thank you for coming.
ID Classifier Metasearch	
Conclusion	

■ のへで

・ロト ・四ト ・ヨト ・ヨト

Thank You

Meta-Evaluating Search Engines

Golbus

Motivation

Diversity

Framework RIC Conditional Rank Correlatio

ID

Classifier Metasearch

Conclusion

Northeastern University

College of Computer and Information Science

Upon recommendation of the President and Faculty and by authority of the Commonwealth of Massachusetts, the Board of Trustees has conferred the degree of

Doctor of Philosophy Computer Beience

apon Peter Bernard Golhus

with all the honors, privileges and responsibilities

appertaining thereunto. Signed and scaled at Boston, Massachusetts, this second day of May in the year two thousand and fourteen.

Joenn Aan

Meta-Evaluating Search Engines

	Thank You
Meta- Evaluating Search Engines Golbus Motivation	
Diversity	
Framework RIC Conditional Rank Correlation	Questions / Comments / Concerns
ID Classifier Metasearch	
Conclusion	

■ のへで

・ロト ・ 日 ト ・ ヨ ト ・ ヨ ト