# Thesaurus-Based Feedback to Support Mixed Search and Browsing Environments

Edgar Meij and Maarten de Rijke

ISLA University of Amsterdam

11th European Conference on Research and Advanced Technology for Digital Libraries (ECDL 2007)





#### **Outline**

- Introduction
  - Motivation
  - Research Questions
  - Language Modeling
- Our Algorithm
  - Overview
  - Determining Thesaurus Terms
  - Estimating a Thesaurus-biased Model
  - Interpolating the Original Query Model
- Results and Discussion
  - Test collection
  - Retrieval Effectiveness
  - Browsing Effectiveness
  - Per-topic Results



#### Outline

- Introduction
  - Motivation
  - Research Questions
  - Language Modeling
- Our Algorithm
  - Overview
  - Determining Thesaurus Terms
  - Estimating a Thesaurus-biased Model
  - Interpolating the Original Query Model
- Results and Discussion
  - Test collection
  - Retrieval Effectiveness
  - Browsing Effectiveness
  - Per-topic Results





- Vocabulary mismatch
  - Not all authors use the same terms
  - Different authors may use different terms for a single concept or may even denote different concepts with the same term
- Solutions
  - Use the cataloging system/controlled vocabulary/thesaurus of the digital library
  - Apply query enrichment/expansion





- Vocabulary mismatch
  - Not all authors use the same terms
  - Different authors may use different terms for a single concept or may even denote different concepts with the same term
- Solutions
  - Use the cataloging system/controlled vocabulary/thesaurus of the digital library
  - Apply query enrichment/expansion





- Information access in a digital library is usually associated with two tasks:
  - Searching
  - Browsing
- Query expansion can be used to improve the search component, whereas a controlled vocabulary may used to enhance the browsing component.
- Aim: How can we combine these?





- Information access in a digital library is usually associated with two tasks:
  - Searching
    - Browsing
- Query expansion can be used to improve the search component, whereas a controlled vocabulary may used to enhance the browsing component.
- Aim: How can we combine these?





- Information access in a digital library is usually associated with two tasks:
  - Searching
    - Browsing
- Query expansion can be used to improve the search component, whereas a controlled vocabulary may used to enhance the browsing component.
- Aim: How can we combine these?





#### **Research Questions**

- How can we use a language modeling framework to generate thesaurus terms, as well as provide query expansion/pseudo-relevance feedback?
- What is the impact of the size of the corpus from which feedback terms are being generated?
- Can our model compete with state-of-the-art IR approaches?
- 4 How can we assess the quality of the thesaurus terms being proposed for browsing?





#### **Outline**

- Introduction
  - Motivation
  - Research Questions
  - Language Modeling
- Our Algorith
  - Overview
  - Determining Thesaurus Terms
  - Estimating a Thesaurus-biased Model
  - Interpolating the Original Query Model
- Results and Discussion
  - Test collection
  - Retrieval Effectiveness
  - Browsing Effectiveness
  - Per-topic Results





## Generative Language Models

• Query-likelihood approach:

$$P(Q|d) \propto P(d) \cdot \prod_{q \in Q} P(q|\theta_d)$$
  
 $\propto \prod_{q \in Q} \frac{c(q,d)}{|d|}$ 

With Dirichlet smoothing:

$$P(Q|d) \propto \prod_{q \in Q} \frac{c(q,d) + \mu P(q|\theta_C)}{|d| + \mu}$$





## Generative Language Models

• Query-likelihood approach:

$$P(Q|d) \propto P(d) \cdot \prod_{q \in Q} P(q|\theta_d)$$
  
  $\propto \prod_{q \in Q} \frac{c(q,d)}{|d|}$ 

With Dirichlet smoothing:

$$P(Q|d) \propto \prod_{q \in Q} \frac{c(q,d) + \mu P(q|\theta_C)}{|d| + \mu}$$





#### Relevance Models

- Generative language modeling assumes that queries are generated from documents
- Relevance modeling assumes both are generated from an unseen source—a relevance model
- A set of documents R is used as a model from which terms are sampled

$$P(w|\hat{\theta}_Q) \propto \sum_{d \in R} P(w|\theta_d) \cdot P(Q|d)$$





#### Relevance Models

- Generative language modeling assumes that queries are generated from documents
- Relevance modeling assumes both are generated from an unseen source—a relevance model
- A set of documents R is used as a model from which terms are sampled

$$P(w|\hat{\theta}_Q) \propto \sum_{d \in R} P(w|\theta_d) \cdot P(Q|d)$$





#### Relevance Models

- Generative language modeling assumes that queries are generated from documents
- Relevance modeling assumes both are generated from an unseen source—a relevance model
- A set of documents R is used as a model from which terms are sampled

$$P(w|\hat{\theta}_Q) \propto \sum_{d \in R} P(w|\theta_d) \cdot P(Q|d)$$





# Final Ranking

- Ranking then comes down to calculating the *distance* between  $P(w|\theta_Q)$  and  $P(w|\theta_d)$  for  $w \in V$
- E.g. using the KL-divergence

$$D_{kl}(\theta_Q||\theta_d) = \sum_{w} P(w|\theta_Q) \cdot \log \frac{P(w|\theta_Q)}{P(w|\theta_d)}$$

## Final Ranking

- Ranking then comes down to calculating the *distance* between  $P(w|\theta_Q)$  and  $P(w|\theta_d)$  for  $w \in V$
- E.g. using the KL-divergence

$$D_{kl}(\theta_Q||\theta_d) = \sum_{w} P(w|\theta_Q) \cdot \log \frac{P(w|\theta_Q)}{P(w|\theta_d)}$$



#### Overview

Determining Thesaurus Terms
Estimating a Thesaurus-biased Model
nterpolating the Original Query Model

#### Outline

- Introduction
  - Motivation
  - Research Questions
  - Language Modeling
- Our Algorithm
  - Overview
  - Determining Thesaurus Terms
  - Estimating a Thesaurus-biased Model
  - Interpolating the Original Query Model
- Results and Discussion
  - Test collection
  - Retrieval Effectiveness
  - Browsing Effectiveness
  - Per-topic Results





#### Overview

Determining Thesaurus Terms
Estimating a Thesaurus-biased Model
Interpolating the Original Query Model

## Our Algorithm: Three Steps

- Determine the thesaurus terms most closely associated with a query
- Search the documents associated with these thesaurus terms, in conjunction with the query, to look for additional terms to describe the query
- Interpolate the query model with the found terms





# **Determining Thesaurus Terms**

 For any given query Q, rank the thesaurus terms m ∈ M according to:

$$P(m|Q) = \frac{P(m)P(Q|m)}{P(Q)}$$
$$= P(m)\sum_{d} P(Q|d)P(d|m)$$

# Estimating a Thesaurus-biased Model

 Then, estimate a thesaurus-biased relevance model by incorporating the top-/ thesaurus terms:

$$P(w|\hat{\theta}_Q) \propto \sum_{d \in R} P(w|\theta_d) \cdot P(Q|d) \cdot P(m_1, \dots, m_l|d)$$

• Assuming the thesaurus terms  $m_1, \ldots, m_l$  to be independent and P(d) to be uniform, we obtain

$$P(w|\hat{\theta}_Q) \propto \sum_{d \in R} P(w|\theta_d) \cdot P(Q|d) \cdot \prod_{i=1,...,l} P(d|m_i) \cdot P(m_i)$$





# Estimating a Thesaurus-biased Model

 Then, estimate a thesaurus-biased relevance model by incorporating the top-/ thesaurus terms:

$$P(w|\hat{\theta}_Q) \propto \sum_{d \in R} P(w|\theta_d) \cdot P(Q|d) \cdot P(m_1, \dots, m_l|d)$$

• Assuming the thesaurus terms  $m_1, \ldots, m_l$  to be independent and P(d) to be uniform, we obtain

$$P(w|\hat{\theta}_Q) \propto \sum_{d \in R} P(w|\theta_d) \cdot P(Q|d) \cdot \prod_{i=1,\dots,J} P(d|m_i) \cdot P(m_i)$$





# Interpolating the Original Query Model

 Finally, the found model is interpolated with the original query using a mixing weight λ to yield the final query model

$$P(w|\theta_Q) = \lambda \cdot \frac{c(w,Q)}{|Q|} + (1-\lambda) \cdot P(w|\hat{\theta}_Q)$$

• When  $\lambda$  is set to 1, a query-likelihood ranking is obtained





# Interpolating the Original Query Model

 Finally, the found model is interpolated with the original query using a mixing weight λ to yield the final query model

$$P(w|\theta_Q) = \lambda \cdot \frac{c(w,Q)}{|Q|} + (1-\lambda) \cdot P(w|\hat{\theta}_Q)$$

• When  $\lambda$  is set to 1, a query-likelihood ranking is obtained





#### Outline

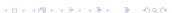
- Introduction
  - Motivation
  - Research Questions
  - Language Modeling
- Our Algorithm
  - Overview
  - Determining Thesaurus Terms
  - Estimating a Thesaurus-biased Model
  - Interpolating the Original Query Model
- Results and Discussion
  - Test collection
  - Retrieval Effectiveness
  - Browsing Effectiveness
  - Per-topic Results





#### TREC Genomics 2006

- Passage retrieval from 160k full-text (biomedical) documents
  - We only look at document-level relevance assessments
  - 28 topics
- PubMed
  - Bibliographic database maintained by the National Library of Medicine (NLM)
  - Over 15M entries, containing author information, abstracts, etc. etc.
- Medical Subject Headings (MeSH)
  - 22,997 hierarchically ordered concepts
  - Trained annotators from the NLM assign one or more
     MeSH terms to every document indexed in PubMed
     UNIVERSITEIT VAN AMSTERDAM



#### TREC Genomics 2006

- Passage retrieval from 160k full-text (biomedical) documents
  - We only look at document-level relevance assessments
  - 28 topics
- PubMed
  - Bibliographic database maintained by the National Library of Medicine (NLM)
  - Over 15M entries, containing author information, abstracts, etc. etc.
- Medical Subject Headings (MeSH)
  - 22,997 hierarchically ordered concepts
  - Trained annotators from the NLM assign one or more
     MeSH terms to every document indexed in PubMed
     UNIVERSITEIT VAN AMSTERDAM
     UNIVERSITEIT VAN AMSTERDAM



#### TREC Genomics 2006

- Passage retrieval from 160k full-text (biomedical) documents
  - We only look at document-level relevance assessments
  - 28 topics
- PubMed
  - Bibliographic database maintained by the National Library of Medicine (NLM)
  - Over 15M entries, containing author information, abstracts, etc. etc.
- Medical Subject Headings (MeSH)
  - 22,997 hierarchically ordered concepts
  - Trained annotators from the NLM assign one or more
     MeSH terms to every document indexed in PubMed
     UNIVERSITEIT VAN AMSTERDAM
     UNIVERSITEIT VAN AMSTERDAM

#### Retrieval Effectiveness

	λ	MAP	P10	
QL	1	0.359	0.45	
RM (collection)	0.10	0.426	+19% 0.48	+7%
RM (PubMed)	0.35	0.425	+18% 0.48	+7%
MM (collection)	0.05	0.424	+18% 0.48	+7%
MM (PubMed)	0.45	0.429	+20% <b>0.49</b>	+9%

Comparison between different query models and a query-likelihood baseline (best scores in boldface.)





#### Retrieval Effectiveness

	λ	MAP	P10	
QL	1	0.359	0.45	
RM (collection)	0.10	0.426	+19% 0.48	+7%
RM (PubMed)	0.35	0.425	+18% 0.48	+7%
MM (collection)	0.05	0.424	+18% 0.48	+7%
MM (PubMed)	0.45	0.429	+20% <b>0.49</b>	+9%

Comparison between different query models and a query-likelihood baseline (best scores in boldface.)



#### Retrieval Effectiveness

	λ	MAP	P10	
QL	1	0.359	0.45	
RM (collection)	0.10	0.426	+19% 0.48	+7%
RM (PubMed)	0.35	0.425	+18% 0.48	+7%
MM (collection)	0.05	0.424	+18% 0.48	+7%
MM (PubMed)	0.45	0.429	+20% <b>0.49</b>	+9%

Comparison between different query models and a query-likelihood baseline (best scores in boldface.)





# **Browsing Effectiveness**

- Concept specificity
- Average distance from each concept to the root of the thesaurus:
  - Collection-based: 4.46
  - PubMed-based: 4.78



# **Browsing Effectiveness**

- TREC 2006 assessors assigned MeSH terms to relevant passages
- Average agreement with the assessors:
  - Collection-based: 2.3/10
  - PubMed-based: 3.0/10
  - Difference is statistically significant (p < 0.05)</li>





# **Browsing Effectiveness**

- TREC 2006 assessors assigned MeSH terms to relevant passages
- Average agreement with the assessors:
  - Collection-based: 2.3/10
  - PubMed-based: 3.0/10
  - Difference is statistically significant (p < 0.05)</li>





#### Summary

- We integrated information inherent in digital libraries into a generative language modeling framework
- Our aim was to facilitate browsing, while maintaining and/or improving retrieval effectiveness
- While readily providing thesaurus terms, our model outperforms state-of-the-art IR methods when estimated on a sufficiently large corpus