# Latent Topic Feedback for Information Retrieval

David Andrzejewski    David Buttler

**Lawrence Livermore National Laboratory**

Center for Applied Scientific Computing
Lawrence Livermore National Laboratory (USA)

August 22, 2011

# BigCo Internal Document Navigation Portal

euro opposition    search

# BigCo Internal Document Navigation Portal

euro opposition    search

## Returned documents

**Hurd in passionate Maastricht defense**
**Financial Times - 14 May 91**

**Small companies may lose in EC deals**
**Financial Times - 14 May 91**

**Russian President Yeltsin invited to G7**
**Financial Times - 24 Mar 92**

- 
- 
-

# BigCo Internal Document Navigation Portal

euro opposition [search]

## Returned documents

**Hurd in passionate Maastricht defense**
**Financial Times - 14 May 91**

**Small companies may lose in EC deals**
**Financial Times - 14 May 91**

**Russian President Yeltsin invited to G7**
**Financial Times - 24 Mar 92**

● 
● 
● 

## Related topics

**debate**
 Tory Euro sceptics
 social chapter, Liberal Democrat
 mps, Labour, bill, Commons

**Emu**
 economic monetary union
 Maastricht treaty, member states
 European, Europe, Community, Emu

● 
● 
●

# Corpus navigation challenges

| Condition | Impaired IR technique |
|---|---|
| Non-expert user | keyword queries |
| Lack of metadata | faceted search |
| Specialized domain | WordNet |
| Small user base | query log mining, relevance feedback |
| Proprietary data | Crowdsourcing |

## Who has these problems?

- Private organizations
- Government agencies

# Corpus navigation challenges

| Condition | Impaired IR technique |
| --- | --- |
| Non-expert user | keyword queries |
| Lack of metadata | faceted search |
| Specialized domain | WordNet |
| Small user base | query log mining, relevance feedback |
| Proprietary data | Crowdsourcing |

## Who has these problems?

- Private organizations
- Government agencies

# Corpus navigation challenges

| Condition | Impaired IR technique |
|---|---|
| Non-expert user | keyword queries |
| Lack of metadata | faceted search |
| Specialized domain | WordNet |
| Small user base | query log mining, relevance feedback |
| Proprietary data | Crowdsourcing |

## Who has these problems?

- Private organizations
- Government agencies

# Corpus navigation challenges

| Condition | Impaired IR technique |
|-----------|----------------------|
| Non-expert user | keyword queries |
| Lack of metadata | faceted search |
| Specialized domain | WordNet |
| Small user base | query log mining, relevance feedback |
| Proprietary data | Crowdsourcing |

## Who has these problems?

- Private organizations
- Government agencies

# Corpus navigation challenges

| Condition | Impaired IR technique |
|---|---|
| Non-expert user | keyword queries |
| Lack of metadata | faceted search |
| Specialized domain | WordNet |
| Small user base | query log mining, relevance feedback |
| Proprietary data | Crowdsourcing |

## Who has these problems?

- Private organizations
- Government agencies

# Corpus navigation challenges

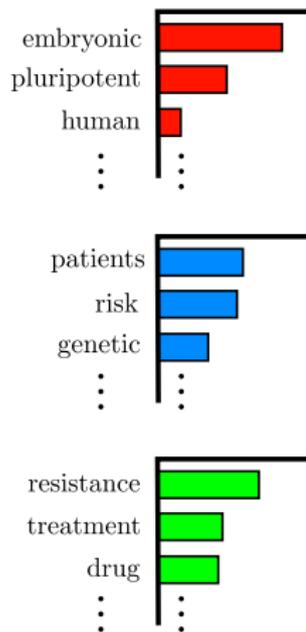| Condition | Impaired IR technique |
|---|---|
| Non-expert user | keyword queries |
| Lack of metadata | faceted search |
| Specialized domain | WordNet |
| Small user base | query log mining, relevance feedback |
| Proprietary data | Crowdsourcing |

Who has these problems?

- Private organizations
- Government agencies

# Corpus navigation challenges

| Condition | Impaired IR technique |
|-----------|----------------------|
| Non-expert user | keyword queries |
| Lack of metadata | faceted search |
| Specialized domain | WordNet |
| Small user base | query log mining, relevance feedback |
| Proprietary data | Crowdsourcing |

## Who has these problems?

- Private organizations
- Government agencies

# Corpus navigation challenges

| Condition | Impaired IR technique |
|-----------|----------------------|
| Non-expert user | keyword queries |
| Lack of metadata | faceted search |
| Specialized domain | WordNet |
| Small user base | query log mining, relevance feedback |
| Proprietary data | Crowdsourcing |

## Who has these problems?

- Private organizations
- Government agencies

Topics $\phi$

Topics $\phi$

Document-topic weights $\theta$

embryonic
pluripotent
human

patients
risk
genetic

resistance
treatment
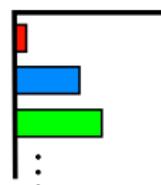drug

Human embryonic stem cell research may benefit patients with genetic risk factors...

Patients at risk for drug-resistant infection...

Topics $\phi$

Document-topic weights $\theta$

embryonic
pluripotent
human

Human embryonic stem cell research may benefit patients with genetic risk factors...

patients
risk
genetic

Patients at risk for drug-resistant infection...

resistance
treatment
drug

Observed **w**    Patients | at | risk | for | drug-resistant
Latent **z**

# How can we exploit latent topics?

- **Implicitly**: language model smoothing (Wei & Croft, SIGIR 2006)
- This approach: **explicit** user feedback on topics
    1. How to show topics?
    2. Which topics to show?
    3. How to use feedback?

# How can we exploit latent topics?

- **Implicitly**: language model smoothing (Wei & Croft, SIGIR 2006)
- This approach: **explicit** user feedback on topics
  1. How to show topics?
  2. Which topics to show?
  3. How to use feedback?

# How can we exploit latent topics?

- **Implicitly**: language model smoothing (Wei & Croft, SIGIR 2006)
- This approach: **explicit** user feedback on topics
  1. How to show topics?
  2. Which topics to show?
  3. How to use feedback?

# How can we exploit latent topics?

- **Implicitly**: language model smoothing (Wei & Croft, SIGIR 2006)
- This approach: **explicit** user feedback on topics
  1. How to show topics?
  2. Which topics to show?
  3. How to use feedback?

# How can we exploit latent topics?

- **Implicitly**: language model smoothing (Wei & Croft, SIGIR 2006)
- This approach: **explicit** user feedback on topics
    1. How to show topics?
    2. Which topics to show?
    3. How to use feedback?

# Question 1 - How to show topics to user?

- "Top N" lists are hard to interpret
- We combine several techniques
    - topic label (Lau et al, COLING 2010)
    - topic *n*-grams (Blei & Lafferty, arXiv 2009)
    - capitalization recovery

| Label | Terms |
| --- | --- |
| Topic 11 | oil, gas, production, exploration sea, north, company, field, energy petroleum, companies |
| Petroleum | state oil company North Sea, natural gas production, exploration, field, energy |

# Question 1 - How to show topics to user?

- "Top N" lists are hard to interpret
- We combine several techniques
  - topic label (Lau et al, COLING 2010)
  - topic *n*-grams (Blei & Lafferty, arXiv 2009)
  - capitalization recovery

| Label | Terms |
|-------|-------|
| Topic 11 | oil, gas, production, exploration<br>sea, north, company, field, energy<br>petroleum, companies |
| Petroleum | state oil company<br>North Sea, natural gas<br>production, exploration, field, energy |

# Question 1 - How to show topics to user?

- "Top N" lists are hard to interpret
- We combine several techniques
  - topic label (Lau et al, COLING 2010)
  - topic *n*-grams (Blei & Lafferty, arXiv 2009)
  - capitalization recovery

| Label | Terms |
|---|---|
| Topic 11 | oil, gas, production, exploration sea, north, company, field, energy petroleum, companies |
| Petroleum | state oil company North Sea, natural gas production, exploration, field, energy |

# Question 1 - How to show topics to user?

- "Top N" lists are hard to interpret
- We combine several techniques
    - topic label (Lau et al, COLING 2010)
    - topic *n*-grams (Blei & Lafferty, arXiv 2009)
    - capitalization recovery

| Label | Terms |
|-------|-------|
| Topic 11 | oil, gas, production, exploration sea, north, company, field, energy petroleum, companies |
| Petroleum | state oil company North Sea, natural gas production, exploration, field, energy |

# Question 1 - How to show topics to user?

- "Top N" lists are hard to interpret
- We combine several techniques
  - topic label (Lau et al, COLING 2010)
  - topic *n*-grams (Blei & Lafferty, arXiv 2009)
  - capitalization recovery

| Label | Terms |
|-------|-------|
| Topic 11 | oil, gas, production, exploration sea, north, company, field, energy petroleum, companies |
| Petroleum | state oil company North Sea, natural gas production, exploration, field, energy |

# Question 2 - Which topics to show?

## Problems

A) Too many topics to present them all ($T > 100$)

B) Incoherent "junk" topics

# Question 2 - Which topics to show?

## Problems

A) Too many topics to present them all ($T > 100$)

B) Incoherent "junk" topics

| | |
|---|---|
| Topic 248 | ve, year, ll, time, don, good, lot, back years, things, make |
| Topic 18 | january, february, december march, month, year, rose feb, sales, fell, increase |

# Problem A - Narrowing down the topics

- Pseudo-relevance feedback → **enriched** topics *E*
- Topic covariance $\Sigma$ → **related** topics *R*
- Top 2 docs, top 2 enriched, top 2 related $\leq$ 12 topics shown

$$E = \bigcup_{d \in D_q} \text{k-}\underset{t}{\text{argmax}}\ \theta_d(t)$$

# Problem A - Narrowing down the topics

- Pseudo-relevance feedback → **enriched** topics $E$
- Topic covariance $\Sigma$ → **related** topics $R$
- Top 2 docs, top 2 enriched, top 2 related $\leq$ 12 topics shown

$$E = \bigcup_{d \in D_q} \text{k-}\underset{t}{\text{argmax}}\ \theta_d(t)$$

$$R = \bigcup_{t \in E} \text{k-}\underset{t' \notin E}{\text{argmax}}\ \Sigma(t, t')$$

## Problem A - Narrowing down the topics

- Pseudo-relevance feedback → **enriched** topics $E$
- Topic covariance $\Sigma$ → **related** topics $R$
- Top 2 docs, top 2 enriched, top 2 related $\leq$ 12 topics shown

$$E = \bigcup_{d \in D_q} \text{k-}\underset{t}{\text{argmax}} \; \theta_d(t)$$

$$R = \bigcup_{t \in E} \text{k-}\underset{t' \notin E}{\text{argmax}} \; \Sigma(t, t')$$

Word co-occurrences in Wikipedia $\rightarrow$ topic PMI score

Incoherent topic

PMI = 0.63

# Problem B - Identifying junk topics
Newman et al (JCDL 2010)

Word co-occurrences in Wikipedia $\rightarrow$ topic PMI score

## Coherent topic
PMI = 3.85

## Incoherent topic
PMI = 0.63

# Problem B - Identifying junk topics
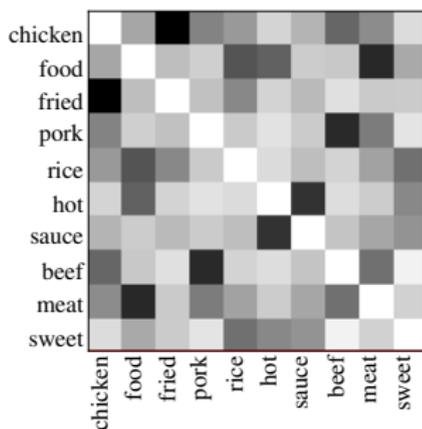
Newman et al (JCDL 2010)

Word co-occurrences in Wikipedia $\rightarrow$ topic PMI score

## Coherent topic

PMI = 3.85



## Incoherent topic

PMI = 0.63

# Problem B - Discarding junk topics

1. Compute PMI scores for each topic $t$
2. Worst PMI scores → **dropped** topics $D$

$$PMI(t) = \frac{1}{k(k-1)} \sum_{(w,w') \in W_t} PMI(w, w')$$

# Problem B - Discarding junk topics

1. Compute PMI scores for each topic $t$
2. Worst PMI scores $\rightarrow$ **dropped** topics $D$

$$PMI(t) = \frac{1}{k(k-1)} \sum_{(w, w') \in W_t} PMI(w, w')$$

$$D = \{t | t \in E \cup R \text{ and } PMI(t) < PMI_{25}\}$$

# Problem B - Discarding junk topics

## Final topics shown

**enriched** and **related**, minus **dropped** $\rightarrow \{E \cup R\} \setminus D$

1. Compute PMI scores for each topic $t$
2. Worst PMI scores $\rightarrow$ **dropped** topics $D$

$$PMI(t) = \frac{1}{k(k-1)} \sum_{(w,w') \in W_t} PMI(w, w')$$

$$D = \{t | t \in E \cup R \text{ and } PMI(t) < PMI_{25}\}$$

# Problem B - Discarding junk topics

## Final topics shown

**enriched** and **related**, minus **dropped** $\rightarrow \{E \cup R\} \setminus D$

1. Compute PMI scores for each topic $t$
2. Worst PMI scores $\rightarrow$ **dropped** topics $D$

$$PMI(t) = \frac{1}{k(k-1)} \sum_{(w,w') \in W_t} PMI(w, w')$$

$$D = \{t | t \in E \cup R \text{ and } PMI(t) < PMI_{25}\}$$

# Question 3 - How to incorporate feedback?

Mechanism should

- **preserve original query intent**
- incorporate the feedback
- "plug and play" with existing search technologies

## Topic-driven query expansion

Weighted combination of

- Original query words $q$
- Top 10 topic words $W_z$

Mechanism should

- preserve original query intent
- incorporate the feedback
- "plug and play" with existing search technologies

## Topic-driven query expansion

Weighted combination of

- Original query words $q$
- Top 10 topic words $W_z$

# Question 3 - How to incorporate feedback?

Mechanism should

- preserve original query intent
- incorporate the feedback
- "plug and play" with existing search technologies

Topic-driven query expansion

Weighted combination of

- Original query words $q$
- Top 10 topic words $W_z$

# Question 3 - How to incorporate feedback?

Mechanism should

- preserve original query intent
- incorporate the feedback
- "plug and play" with existing search technologies

## Topic-driven query expansion

Weighted combination of

- Original query words $q$
- Top 10 topic words $W_z$

# Question 3 - How to incorporate feedback?

Mechanism should

- preserve original query intent
- incorporate the feedback
- "plug and play" with existing search technologies

## Topic-driven query expansion

Weighted combination of

- Original query words $q$
- Top 10 topic words $W_z$

# Question 3 - How to incorporate feedback?

Mechanism should

- preserve original query intent
- incorporate the feedback
- "plug and play" with existing search technologies

## Topic-driven query expansion

Weighted combination of

- Original query words $q$
- Top 10 topic words $W_z$

# Example TREC query

- Corpus: 210K news articles (Financial Times, 1992-1994)
- Query: "euro opposition"
  (political opposition to the € currency union)
- Ground truth: 98 articles judged **relevant**

# Example TREC query

- Corpus: 210K news articles (Financial Times, 1992-1994)
- Query: "euro opposition"
  (political opposition to the € currency union)
- Ground truth: 98 articles judged **relevant**

# Example TREC query

- Corpus: 210K news articles (Financial Times, 1992-1994)
- Query: "euro opposition"
  (political opposition to the € currency union)
- Ground truth: 98 articles judged **relevant**

## "euro opposition" topics

| Label | Terms | PMI percentile |
|-------|-------|----------------|
| debate | Tory Euro sceptics<br>social chapter, Liberal Democrat<br>mps, Labour, bill, Commons | 47 |
| business | PERSONAL FILE Born<br>years ago, past years<br>man, time, job, career | **2** |
| Emu | economic monetary union<br>Maastricht treaty, member states<br>European, Europe, Community, Emu | 63 |
| George | President George Bush, White House<br>Mr Clinton, administration<br>Democratic, Republican, Washington | 60 |

# "euro opposition" topics

| Label | Terms | PMI percentile |
|-------|-------|----------------|
| debate | Tory Euro sceptics<br>social chapter, Liberal Democrat<br>mps, Labour, bill, Commons | 47 |
| ~~business~~ | PERSONAL FILE Born<br>years ago, past years<br>man, time, job, career | **2** |
| Emu | economic monetary union<br>Maastricht treaty, member states<br>European, Europe, Community, Emu | 63 |
| George | President George Bush, White House<br>Mr Clinton, administration<br>Democratic, Republican, Washington | 60 |

# "euro opposition" topics

| Label | Terms | PMI percentile |
|-------|-------|----------------|
| debate | Tory Euro sceptics<br>social chapter, Liberal Democrat<br>mps, Labour, bill, Commons | 47 |
| ~~business~~ | PERSONAL FILE Born<br>years ago, past years<br>man, time, job, career | **2** |
| **Emu** | economic monetary union<br>Maastricht treaty, member states<br>European, Europe, Community, Emu | 63 |
| George | President George Bush, White House<br>Mr Clinton, administration<br>Democratic, Republican, Washington | 60 |

## "euro opposition" topics

| Label | Terms | PMI percentile |
|-------|-------|----------------|
| debate | Tory Euro sceptics<br>social chapter, Liberal Democrat<br>mps, Labour, bill, Commons | 47 |
| ~~business~~ | PERSONAL FILE Born<br>years ago, past years<br>man, time, job, career | **2** |
| **Emu** | economic monetary union<br>Maastricht treaty, member states<br>European, Europe, Community, Emu | 63 |
| George | President George Bush, White House<br>Mr Clinton, administration<br>Democratic, Republican, Washington | 60 |

## "euro opposition" topics

| Label | Terms | PMI percentile |
|---|---|---|
| debate | Tory Euro sceptics<br>social chapter, Liberal Democrat<br>mps, Labour, bill, Commons | 47 |
| ~~business~~ | PERSONAL FILE Born<br>years ago, past years<br>man, time, job, career | **2** |
| **Emu** | economic monetary union<br>Maastricht treaty, member states<br>European, Europe, Community, Emu | 63 |
| George | President George Bush, White House<br>Mr Clinton, administration<br>Democratic, Republican, Washington | 60 |

# "Emu" topic feedback

Indri weighted query operator                    Original query

```
#weight(0.375 euro, 0.375 opposition,
        0.031 European, ..., 0.015 Emu)
```

Topic expansion                    ROC curve (true/false positive rates)

| Measure | Gain |
|---------|------|
| NDCG15  | +0.22 |
| NDCG    | +0.07 |
| MAP     | +0.02 |

# "Emu" topic feedback

Indri weighted query operator — Original query

```
#weight (0.375 euro, 0.375 opposition,
         0.031 European, ..., 0.015 Emu)
```
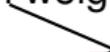
Topic expansion — ROC curve (true/false positive rates)

| Measure | Gain |
|---------|------|
| NDCG15  | +0.22 |
| NDCG    | +0.07 |
| MAP     | +0.02 |

# "Emu" topic feedback

Indri weighted query operator
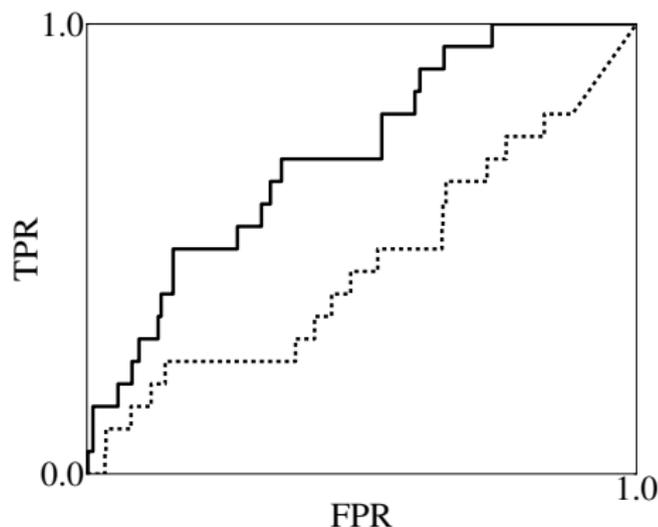
Original query

```
#weight( 0.375 euro, 0.375 opposition,
         0.031 European, ..., 0.015 Emu)
```

Topic expansion                    ROC curve (true/false positive rates)

| Measure | Gain |
|---------|------|
| NDCG15  | +0.22 |
| NDCG    | +0.07 |
| MAP     | +0.02 |

# "Emu" topic feedback

```
#weight(0.375 euro, 0.375 opposition,

        0.031 European, ..., 0.015 Emu )
```

Topic expansion ⟵                    ROC curve (true/false positive rates)

| Measure | Gain |
|---------|------|
| NDCG15  | +0.22 |
| NDCG    | +0.07 |
| MAP     | +0.02 |

# "Emu" topic feedback

```
#weight(0.375 euro, 0.375 opposition,
        0.031 European, ..., 0.015 Emu)
```

Topic expansion                    ROC curve (true/false positive rates)

| Measure | Gain |
|---------|------|
| NDCG15  | $+0.22$ |
| NDCG    | $+0.07$ |
| MAP     | $+0.02$ |

# Experimental results

- TREC datasets
  - 6 newswire corpora, 814K documents total
  - Learn $T = 500$ topics per corpus
  - 850 queries total (some overlap)
- Assume user will select "right" topic (if presented)
- Summary ($h =$ a *helpful* topic exists, $s =$ we show it to the user)
  - Avg number of topics shown $= 7.76$
  - $P(h) \approx 40\%, P(s|h) \approx 40\%$

# Experimental results

- TREC datasets
  - 6 newswire corpora, 814K documents total
  - Learn $T = 500$ topics per corpus
  - 850 queries total (some overlap)
- Assume user will select "right" topic (if presented)
- Summary ($h$ = a *helpful* topic exists, $s$ = we show it to the user)
  - Avg number of topics shown = 7.76
  - $P(h) \approx 40\%, P(s|h) \approx 40\%$
  - 
  - 
  - 
  -

# Experimental results

- TREC datasets
  - 6 newswire corpora, 814K documents total
  - Learn $T = 500$ topics per corpus
  - 850 queries total (some overlap)
- Assume user will select "right" topic (if presented)
- Summary ($h = $ a *helpful* topic exists, $s = $ we show it to the user)
  - Avg number of topics shown = 7.76
  - $P(h) \approx 40\%$, $P(s|h) \approx 40\%$
  - 
  - 
  - 
  -

# Experimental results

- TREC datasets
  - 6 newswire corpora, 814K documents total
  - Learn $T = 500$ topics per corpus
  - 850 queries total (some overlap)
- Assume user will select "right" topic (if presented)
- Summary ($h =$ a *helpful* topic exists, $s =$ we show it to the user)
  - Avg number of topics shown $= 7.76$
  - $P(h) \approx 40\%, P(s|h) \approx 40\%$
  - 
  - 
  - 
  -

# Experimental results

- TREC datasets
    - 6 newswire corpora, 814K documents total
    - Learn $T = 500$ topics per corpus
    - 850 queries total (some overlap)
- Assume user will select "right" topic (if presented)
- Summary ($h$ = a *helpful* topic exists, $s$ = we show it to the user)
    - Avg number of topics shown = 7.76
    - $P(h) \approx 40\%$, $P(s|h) \approx 40\%$

# Experimental results

- TREC datasets
    - 6 newswire corpora, 814K documents total
    - Learn $T = 500$ topics per corpus
    - 850 queries total (some overlap)
- Assume user will select "right" topic (if presented)
- Summary ($h =$ a *helpful* topic exists, $s =$ we show it to the user)
    - Avg number of topics shown $= 7.76$
    - $P(h) \approx 40\%$, $P(s|h) \approx 40\% \rightarrow P(h \wedge s) = 15.6\%$
    - Adding **related** topics helps
      (else $P(h \wedge s) = 10.9\%$, avg shown $= 2.70$)
    - Discarding **dropped** does not hurt
      (else $P(h \wedge s) = 16.8\%$, avg shown $= 9.79$)

- TREC datasets
    - 6 newswire corpora, 814K documents total
    - Learn $T = 500$ topics per corpus
    - 850 queries total (some overlap)
- Assume user will select "right" topic (if presented)
- Summary ($h =$ a *helpful* topic exists, $s =$ we show it to the user)
    - Avg number of topics shown $= 7.76$
    - $P(h) \approx 40\%$, $P(s|h) \approx 40\% \rightarrow P(h \wedge s) = 15.6\%$
    - Adding **related** topics helps
      (else $P(h \wedge s) = 10.9\%$, avg shown $= 2.70$)
    - Discarding **dropped** does not hurt
      (else $P(h \wedge s) = 16.8\%$, avg shown $= 9.79$)

# Experimental results

- TREC datasets
  - 6 newswire corpora, 814K documents total
  - Learn $T = 500$ topics per corpus
  - 850 queries total (some overlap)
- Assume user will select "right" topic (if presented)
- Summary ($h = $ a *helpful* topic exists, $s = $ we show it to the user)
  - Avg number of topics shown $= 7.76$
  - $P(h) \approx 40\%, P(s|h) \approx 40\% \rightarrow P(h \wedge s) = 15.6\%$
  - Adding **related** topics helps
    (else $P(h \wedge s) = 10.9\%$, avg shown $= 2.70$)
  - Discarding **dropped** does not hurt
    (else $P(h \wedge s) = 16.8\%$, avg shown $= 9.79$)

# Experimental results

- TREC datasets
  - 6 newswire corpora, 814K documents total
  - Learn $T = 500$ topics per corpus
  - 850 queries total (some overlap)
- Assume user will select "right" topic (if presented)
- Summary ($h =$ a *helpful* topic exists, $s =$ we show it to the user)
  - Avg number of topics shown $= 7.76$
  - $P(h) \approx 40\%, P(s|h) \approx 40\% \rightarrow P(h \wedge s) = 15.6\%$
  - Adding **related** topics helps
    (else $P(h \wedge s) = 10.9\%$, avg shown $= 2.70$)
  - Discarding **dropped** does not hurt
    (else $P(h \wedge s) = 16.8\%$, avg shown $= 9.79$)

# Experimental results

- TREC datasets
    - 6 newswire corpora, 814K documents total
    - Learn $T = 500$ topics per corpus
    - 850 queries total (some overlap)
- Assume user will select "right" topic (if presented)
- Summary ($h$ = a *helpful* topic exists, $s$ = we show it to the user)
    - Avg number of topics shown = 7.76
    - $P(h) \approx 40\%$, $P(s|h) \approx 40\% \rightarrow P(h \wedge s) = 15.6\%$
    - Adding **related** topics helps
      (else $P(h \wedge s) = 10.9\%$, avg shown = 2.70)
    - Discarding **dropped** does not hurt
      (else $P(h \wedge s) = 16.8\%$, avg shown = 9.79)

# Experimental results

- TREC datasets
  - 6 newswire corpora, 814K documents total
  - Learn $T = 500$ topics per corpus
  - 850 queries total (some overlap)
- Assume user will select "right" topic (if presented)
- Summary ($h$ = a *helpful* topic exists, $s$ = we show it to the user)
  - Avg number of topics shown = 7.76
  - $P(h) \approx 40\%, P(s|h) \approx 40\% \rightarrow P(h \wedge s) = 15.6\%$
  - Adding **related** topics helps
    (else $P(h \wedge s) = 10.9\%$, avg shown = 2.70)
  - Discarding **dropped** does not hurt
    (else $P(h \wedge s) = 16.8\%$, avg shown = 9.79)

# Experimental results

- TREC datasets
  - 6 newswire corpora, 814K documents total
  - Learn $T = 500$ topics per corpus
  - 850 queries total (some overlap)
- Assume user will select "right" topic (if presented)
- Summary ($h$ = a *helpful* topic exists, $s$ = we show it to the user)
  - Avg number of topics shown = 7.76
  - $P(h) \approx 40\%, P(s|h) \approx 40\% \rightarrow P(h \wedge s) = 15.6\%$
  - Adding **related** topics helps
    (else $P(h \wedge s) = 10.9\%$, avg shown = 2.70)
  - Discarding **dropped** does not hurt
    (else $P(h \wedge s) = 16.8\%$, avg shown = 9.79)

# Important idea

Even when topics do *not* improve NDCG and friends . . .
they still may be useful/informative.

# Important idea

Even when topics do *not* improve NDCG and friends . . .

they still may be useful/informative.

# Conclusions and Future Work

- Conclusions
  - Explicit topic feedback can improve relevance
  - Selection approach can find relevant topics
- Future work
  - Better topics? (fancier topic models / user guidance)
  - Better topic selection? (user modeling, learning to rank)
  - Validate assumptions and presentation strategy (user study)
  - Compare / combine with **implicit** topic usage

# Conclusions and Future Work

- Conclusions
  - Explicit topic feedback can improve relevance
  - Selection approach can find relevant topics
- Future work
  - Better topics? (fancier topic models / user guidance)
  - Better topic selection? (user modeling, learning to rank)
  - Validate assumptions and presentation strategy (user study)
  - Compare / combine with **implicit** topic usage

# Conclusions and Future Work

- Conclusions
  - Explicit topic feedback can improve relevance
  - Selection approach can find relevant topics
- Future work
  - Better topics? (fancier topic models / user guidance)
  - Better topic selection? (user modeling, learning to rank)
  - Validate assumptions and presentation strategy (user study)
  - Compare / combine with **implicit** topic usage

# Conclusions and Future Work

- Conclusions
  - Explicit topic feedback can improve relevance
  - Selection approach can find relevant topics
- Future work
  - Better topics? (fancier topic models / user guidance)
  - Better topic selection? (user modeling, learning to rank)
  - Validate assumptions and presentation strategy (user study)
  - Compare / combine with **implicit** topic usage

# Conclusions and Future Work

- Conclusions
  - Explicit topic feedback can improve relevance
  - Selection approach can find relevant topics
- Future work
  - Better topics? (fancier topic models / user guidance)
  - Better topic selection? (user modeling, learning to rank)
  - Validate assumptions and presentation strategy (user study)
  - Compare / combine with **implicit** topic usage

# Conclusions and Future Work

- Conclusions
  - Explicit topic feedback can improve relevance
  - Selection approach can find relevant topics
- Future work
  - Better topics? (fancier topic models / user guidance)
  - Better topic selection? (user modeling, learning to rank)
  - Validate assumptions and presentation strategy (user study)
  - Compare / combine with **implicit** topic usage

# Conclusions and Future Work

- Conclusions
  - Explicit topic feedback can improve relevance
  - Selection approach can find relevant topics
- Future work
  - Better topics? (fancier topic models / user guidance)
  - Better topic selection? (user modeling, learning to rank)
  - Validate assumptions and presentation strategy (user study)
  - Compare / combine with **implicit** topic usage

# Conclusions and Future Work

- Conclusions
  - Explicit topic feedback can improve relevance
  - Selection approach can find relevant topics
- Future work
  - Better topics? (fancier topic models / user guidance)
  - Better topic selection? (user modeling, learning to rank)
  - Validate assumptions and presentation strategy (user study)
  - Compare / combine with **implicit** topic usage

# Demo Web Interface

WED AUG 17 11:49:51 PDT 2011

highway funds | Submit

## GOVERNOR'S BUDGET ASKS $47.8 BILLION; NOT ADEQUATE, HE SAYS, CITING RESTRICTIONS

AUTHOR: AUTHORNAME | INSTITUTION: AFFILIATIONS | PUBLICATION: PUBLISHED BY | DATE: PUBLICATION DATE | ABSTRACT »

## ORANGE COUNTY FOCUS: BREA; CITY STAFF LAUDED AS COUNCIL OKS BUDGET

AUTHOR: AUTHORNAME | INSTITUTION: AFFILIATIONS | PUBLICATION: PUBLISHED BY | DATE: PUBLICATION DATE | ABSTRACT »

## Topics

Refine Query

Your Query: highway funds

q=highway+funds&defType=dismax

Executing query now...

- Topic: 408
  - car pool lanes
  - car pool, san diego
  - traffic, freeway, road, highway
- Topic: 76
  - estimated cost million

## Acknowledgments

- Web demo: Kevin Lawrence (Florida A&M)
- This work was performed under the auspices of the U.S. Department of Energy by Lawrence Livermore National Laboratory under Contract DE-AC52-07NA27344. LLNL-PRES-491932

# Demo Web Interface

WED AUG 17 11:49:51 PDT 2011

Refine Query

Your Query: highway funds

q=highway+funds&defType=dismax

Executing query now...

[ highway funds ] **Submit**

## GOVERNOR'S BUDGET ASKS $47.8 BILLION; NOT ADEQUATE, HE SAYS, CITING RESTRICTIONS

AUTHOR: AUTHORNAME | INSTITUTION: AFFILIATIONS | PUBLICATION: PUBLISHED BY | DATE: PUBLICATION DATE |
ABSTRACT »

## ORANGE COUNTY FOCUS: BREA; CITY STAFF LAUDED AS COUNCIL OKS BUDGET

AUTHOR: AUTHORNAME | INSTITUTION: AFFILIATIONS | PUBLICATION: PUBLISHED BY | DATE: PUBLICATION DATE |
ABSTRACT »

○ Topic: 408
- car pool lanes
- car pool, san diego
- traffic, freeway, road, highway

○ Topic: 76
- estimated cost million

## Acknowledgments

- Web demo: Kevin Lawrence (Florida A&M)
- This work was performed under the auspices of the U.S. Department of Energy by Lawrence Livermore National Laboratory under Contract DE-AC52-07NA27344. LLNL-PRES-491932

# Predicting relevant topics

Can we *predict* which topics will improve relevance?

- Short answer: no (well, I couldn't get it to work. . .)
- Linear / logistic regression

| Feature | Interpretation |
| --- | --- |
| $PMI(t)$ | topic quality |
| $Entropy(P(d|t))$ | document-concentration of topic |
| $\log(P(q|t))$ | query probability under the topic |
| $\log(\sum_{d \in D_q} \theta_d(t))$ | topic probability across top documents |

Missed helpful topics: "far" from top baseline documents

# Predicting relevant topics

Can we *predict* which topics will improve relevance?

- Short answer: no (well, I couldn't get it to work. . .)
- Linear / logistic regression

| Feature | Interpretation |
|---|---|
| $PMI(t)$ | topic quality |
| $Entropy(P(d\|t))$ | document-concentration of topic |
| $\log(P(q\|t))$ | query probability under the topic |
| $\log(\sum_{d \in D_q} \theta_d(t))$ | topic probability across top documents |

Missed helpful topics: "far" from top baseline documents

# Predicting relevant topics

Can we *predict* which topics will improve relevance?

- Short answer: no (well, I couldn't get it to work...)
- Linear / logistic regression

| Feature | Interpretation |
| --- | --- |
| $PMI(t)$ | topic quality |
| $Entropy(P(d|t))$ | document-concentration of topic |
| $\log(P(q|t))$ | query probability under the topic |
| $\log(\sum_{d \in D_q} \theta_d(t))$ | topic probability across top documents |

Missed helpful topics: "far" from top baseline documents

# Predicting relevant topics

Can we *predict* which topics will improve relevance?

- Short answer: no (well, I couldn't get it to work. . .)
- Linear / logistic regression

| Feature | Interpretation |
|---------|----------------|
| $PMI(t)$ | topic quality |
| $Entropy(P(d\|t))$ | document-concentration of topic |
| $\log(P(q\|t))$ | query probability under the topic |
| $\log(\sum_{d \in D_q} \theta_d(t))$ | topic probability across top documents |

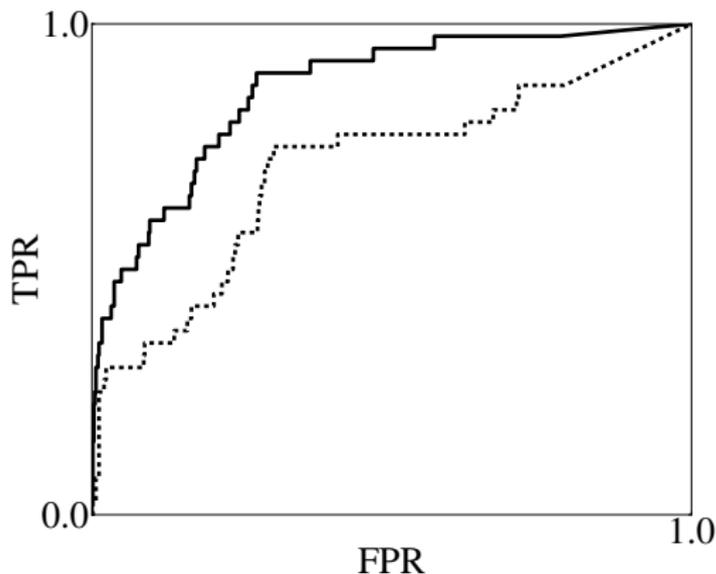Missed helpful topics: "far" from top baseline documents

# Negative feedback

- Could also allow user to mark topic as **not** relevant
- Use Indri `#not` operator to form new query
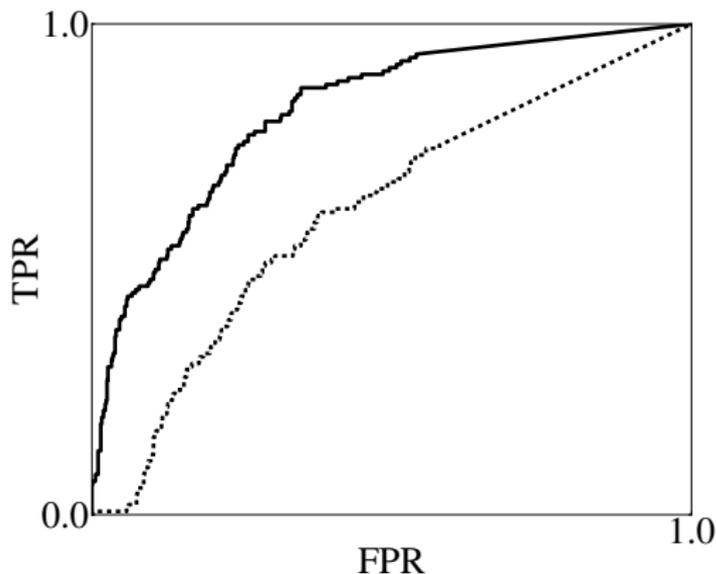- Intuitively appealing, but did not seem to help in experiments...

# Negative feedback

- Could also allow user to mark topic as **not** relevant
- Use Indri `#not` operator to form new query
- Intuitively appealing, but did not seem to help in experiments...

# Negative feedback

- Could also allow user to mark topic as **not** relevant
- Use Indri `#not` operator to form new query
- Intuitively appealing, but did not seem to help in experiments...

# "law enforcement dogs"

| Label | Terms |
|-------|-------|
| heroin | seized kg cocaine, drug traffickers, kg heroin, police, arrested, drugs, marijuana |

## "King Hussein, peace"

| Label | Terms |
|-------|-------|
| Amman | Majesty King Husayn, al Aqabah, peace process, Jordan, Jordanian, Amman, Arab |

## "bank failures"

| Label | Terms |
| --- | --- |
| FDIC | Federal Deposit Insurance, William Seidman, Insurance Corp, banks, bank, FDIC, banking |

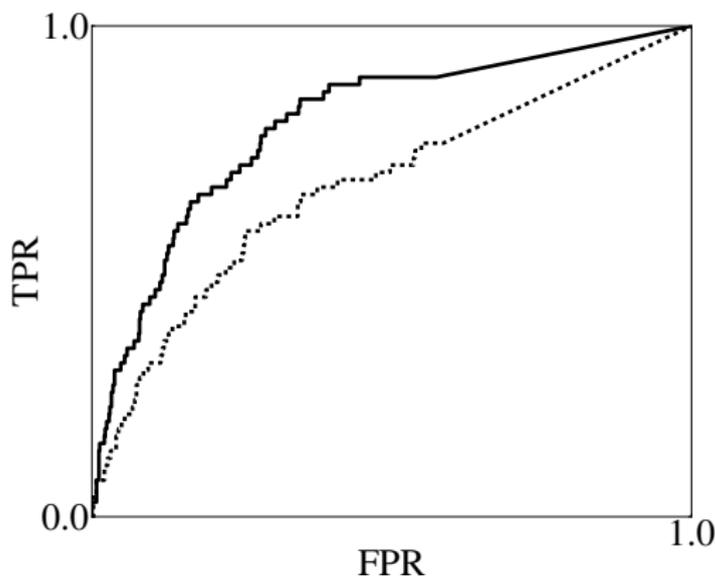# "US-USSR Arms Control Agreements"

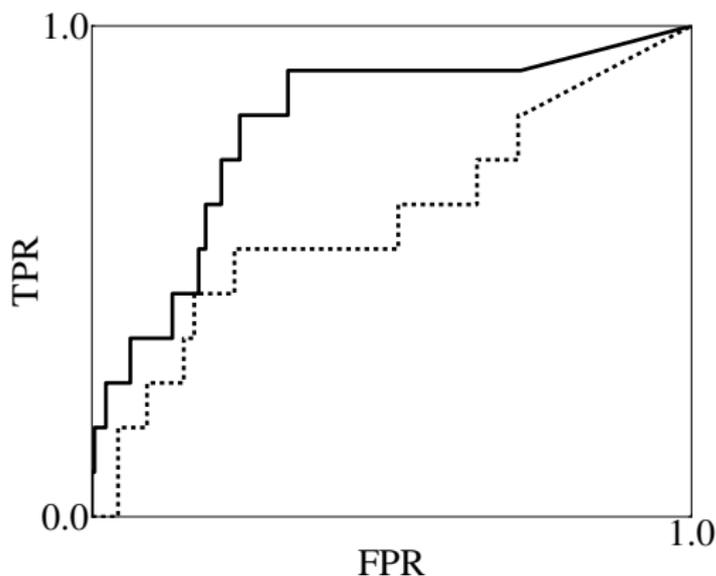| Label | Terms |
|-------|-------|
| missile | Strategic Defense Initiative, United States, arms control, treaty, nuclear, missiles, range |

## "Possible Contributions of Gene Mapping to Medicine"

| Label | Terms |
|-------|-------|
| called | British journal Nature, immune system, genetically engineered, cells, research, researchers, scientists |

## "New Space Satellite Applications"

| Label | Terms |
|-------|-------|
| communications | European Space Agency, Air Force, Cape Canaveral, satellite, launch, rocket, satellites |

... governmental strategy of attracting foreign direct investment,...

...governmental strategy of attracting foreign direct investment,...

...governmental strategy of attracting foreign direct investment,...

...governmental strategy of attracting foreign direct investment ,...