

# Topical Sequence Profiling

---

Tim Gollub\*   Nedim Lipka†   Eunyee Koh†   Erdan Genc\*   Benno Stein\*

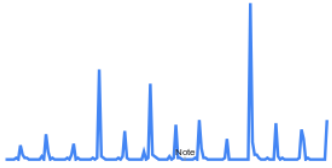
TIR@DEXA   5. Sept. 2016

\* Webis Group  
Bauhaus-Universität Weimar  
[www.webis.de](http://www.webis.de)

† Big Data Experience Lab  
Adobe Systems  
[www.research.adobe.com](http://www.research.adobe.com)

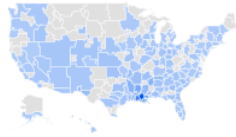
Featured insights

MTV Video Music Awards 2016



Search interest in the VMAs

Louisiana floods



Search interest by metro area, past week

Stories trending now

1 Davao, Philippines, Rodrigo Duterte



2 Buffalo Bills, Manny Lawson, Jerome Felton



3 Food and Drug Administration, Antibacterial soap, Antibio...



4 Jon Polito, Coen brothers, The Big Lebowski



5 Tennessee Titans, Justin Hunter, Bishop Sankey, NFL



6 Los Angeles Dodgers, Vin Scully, KTLA



7 Banksy, Robert Del Naja, Massive Attack

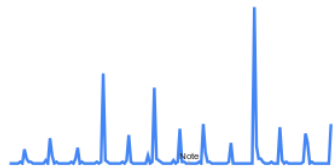


8 Anthony Rizzo, Chicago Cubs, Cancer



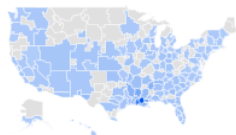
Featured insights

MTV Video Music Awards 2016



Search interest in the VMAs

Louisiana floods



Search interest by metro area, past week

Stories trending now

1 Davao, Philippines, Rodrigo Duterte



2 Buffalo Bills, Manny Lawson, Jerome Felton



3 Food and Drug Administration, Antibacterial soap, Antibio...



4 Jon Polito, Coen brothers, The Big Lebowski



5 Tennessee Titans, Justin Hunter, Bishop Sankey, NFL



6 Los Angeles Dodgers, Vin Scully, KTLA



7 Banksy, Robert Del Naja, Massive Attack



8 Anthony Rizzo, Chicago Cubs, Cancer



WHAT DID THE WORLD SEARCH FOR IN 2012?  
WATCH THE YEAR IN REVIEW

See search trends from around the world.

Select a country

2012 Search Trends

United States

Trending

Searches

- Whitney Houston
- Hurricane Sandy
- Election 2012
- Hunger Games
- Jeremy Lin
- Olympics 2012
- Amanda Todd
- Gangnam Style
- Michael Clarke Duncan
- KONY 2012

Trending

People

- Whitney Houston
- Jeremy Lin
- Amanda Todd
- Michael Clark Duncan
- Kate Middleton
- One Direction
- Morgan Freeman
- Peyton Manning
- Joe Paterno
- Paul Ryan

Trending

Events of 2012

- Hurricane Sandy
- Presidential Election
- Super Bowl
- Olympics
- UEFA Euro 2012
- KONY Movement

Trending

How to...

- How to love
- How to rock
- How to vote
- How to install
- How to hate
- How to archer

# Topical Sequence Profiling

Resource

Stream

Collection

Register

Query

Social Media

Articles, Papers

Books

# Topical Sequence Profiling

## Resource

### Stream

### Collection

Query

Google Trends

Google Zeitgeist

Social Media

Articles, Papers

Books

Register

Time range:

incomplete

complete

Analysis:

differential

covering

# Topical Sequence Profiling

## Resource

### Stream

### Collection

Register

Query

Google Trends

Google Zeitgeist

Social Media

Retweet statistics in  
Twitter

Wikipedia reverts  
analysis

Articles, Papers

Provenance analytics

Topical Sequence  
Profiling

Books

–

Discourse analysis

Time range:

incomplete

complete

Analysis:

differential

covering

# Topical Sequence Profiling

## Problem Statement

Given a sequence  $\mathcal{D} = \{D_1, \dots, D_n\}$  of text collections,

such as

- ❑ a series of the annual proceedings of a conference,
- ❑ a news feed of a certain period, or
- ❑ the social media mentions of an entity over time frames,

provide (statistical) insights about its content.

# Topical Sequence Profiling

## Problem Statement (continued)

A topic embedding  $\mathbf{T}$  of  $\mathcal{D} = \{D_1, \dots, D_n\}$ :

$$\mathbf{T} = \begin{array}{c} t_1 \\ \vdots \\ t_k \end{array} \begin{array}{c} \left[ \begin{array}{ccc} D_1 & \dots & D_n \\ \mathbf{T}_{11} & & \mathbf{T}_{1n} \\ & \dots & \\ \mathbf{T}_{k1} & & \mathbf{T}_{kn} \end{array} \right] \end{array}$$

---

Coverage  $c_1 \quad \dots \quad c_n$



# Topical Sequence Profiling

## Problem Statement (continued)

A topical embedding  $\mathbf{T}$  of  $\mathcal{D} = \{D_1, \dots, D_n\}$ :

		$D_1$	$\dots$	$D_n$
$\mathbf{T} =$	Cluster analysis	4%		40%
	$\vdots$		$\dots$	
	Retrieval model	7%		$\mathbf{T}_{kn}$
	Coverage	72%	$\dots$	$C_n$

# Topical Sequence Profiling

## Problem Statement (continued)

A topic embedding  $\mathbf{T}$  of  $\mathcal{D} = \{D_1, \dots, D_n\}$ :

$$\mathbf{T} = \begin{array}{c} \text{Cluster analysis} \\ \vdots \\ \text{Retrieval model} \end{array} \left[ \begin{array}{ccc} D_1 & \dots & D_n \\ 4\% & & 40\% \\ & \dots & \\ 7\% & & \mathbf{T}_{kn} \end{array} \right]$$

---

Coverage                      72%                      ...                       $c_n$

A topic embedding  $\mathbf{T}$  characterizes  $\mathcal{D} = \{D_1, \dots, D_n\}$  if:

1.  $\mathbf{T}$  is *representative* for each of the text collections  $D \in \mathcal{D}$ .
2.  $\mathbf{T}$  highlights *informative topic developments* within  $\mathcal{D}$ .
3.  $\mathbf{T}$  is *small* enough to be surveyed quickly.

# Topical Sequence Profiling

## Problem Statement (continued)

A topic embedding  $\mathbf{T}$  of  $\mathcal{D} = \{D_1, \dots, D_n\}$ :

$$\mathbf{T} = \begin{array}{l} \text{Cluster analysis} \\ \vdots \\ \text{Retrieval model} \end{array} \left[ \begin{array}{ccc} D_1 & \dots & D_n \\ 4\% & & 40\% \\ & \dots & \\ 7\% & & \mathbf{T}_{kn} \end{array} \right]$$

---

Coverage                      72%                      ...                       $c_n$

### 1. *Representativeness.*

A topic embedding  $\mathbf{T} = \{t_1, \dots, t_k\}$  is representative if for every collection  $D \in \mathcal{D}$  the percentage of documents that address at least one of the topics is above a predefined threshold  $c \in [0; 1]$ .

# Topical Sequence Profiling

## Problem Statement (continued)

A topic embedding  $\mathbf{T}$  of  $\mathcal{D} = \{D_1, \dots, D_n\}$ :

$\mathbf{T} =$	Cluster analysis	$\vdots$	Retrieval model	<table style="border-collapse: collapse; width: 100%; text-align: center;"> <tr> <td style="padding: 5px;"><math>D_1</math></td> <td style="padding: 5px;"><math>\dots</math></td> <td style="padding: 5px;"><math>D_n</math></td> </tr> <tr> <td style="padding: 5px; color: green;">4%</td> <td style="padding: 5px;"><math>\dots</math></td> <td style="padding: 5px; color: green;">40%</td> </tr> <tr> <td style="padding: 5px;">7%</td> <td style="padding: 5px;"><math>\dots</math></td> <td style="padding: 5px;"><math>\mathbf{T}_{kn}</math></td> </tr> </table>	$D_1$	$\dots$	$D_n$	4%	$\dots$	40%	7%	$\dots$	$\mathbf{T}_{kn}$
$D_1$	$\dots$	$D_n$											
4%	$\dots$	40%											
7%	$\dots$	$\mathbf{T}_{kn}$											
<table style="width: 100%; text-align: center;"> <tr> <td style="padding: 5px;">Coverage</td> <td style="padding: 5px;">72%</td> <td style="padding: 5px;"><math>\dots</math></td> <td style="padding: 5px;"><math>C_n</math></td> </tr> </table>				Coverage	72%	$\dots$	$C_n$						
Coverage	72%	$\dots$	$C_n$										

### 2. Informativeness.

The informativeness of a topic’s development is assessed by the variance of the topic distribution. The higher the variance, the more informative is its development. The mean topic variance is used to assess the informativeness of the whole embedding  $\mathbf{T}$ .

$$\frac{1}{k} \sum_{i=1}^k \text{Var}(\mathbf{T}_{i:})$$

# Topical Sequence Profiling

## Problem Statement (continued)

A topical embedding  $\mathbf{T}$  of  $\mathcal{D} = \{D_1, \dots, D_n\}$ :

$\mathbf{T} =$	<div style="display: flex; flex-direction: column; align-items: center;"> <div style="color: green; margin-bottom: 5px;">Cluster analysis</div> <div style="margin-bottom: 5px;"><math>\vdots</math></div> <div style="margin-bottom: 5px;">Retrieval model</div> </div>	[	<table style="width: 100%; border-collapse: collapse;"> <tr> <td style="text-align: center; padding: 5px;"><math>D_1</math></td> <td style="text-align: center; padding: 5px;"><math>\dots</math></td> <td style="text-align: center; padding: 5px;"><math>D_n</math></td> </tr> <tr> <td style="text-align: center; color: green; padding: 5px;">4%</td> <td style="text-align: center; padding: 5px;"><math>\dots</math></td> <td style="text-align: center; color: green; padding: 5px;">40%</td> </tr> <tr> <td style="text-align: center; padding: 5px;">7%</td> <td style="text-align: center; padding: 5px;"><math>\dots</math></td> <td style="text-align: center; padding: 5px;"><math>\mathbf{T}_{kn}</math></td> </tr> </table>	$D_1$	$\dots$	$D_n$	4%	$\dots$	40%	7%	$\dots$	$\mathbf{T}_{kn}$	]
$D_1$	$\dots$	$D_n$											
4%	$\dots$	40%											
7%	$\dots$	$\mathbf{T}_{kn}$											
Coverage		<table style="width: 100%; border-collapse: collapse;"> <tr> <td style="text-align: center; padding: 5px;">72%</td> <td style="text-align: center; padding: 5px;"><math>\dots</math></td> <td style="text-align: center; color: orange; padding: 5px;"><math>C_n</math></td> </tr> </table>			72%	$\dots$	$C_n$						
72%	$\dots$	$C_n$											

### 3. Minimality.

A topical embedding  $\mathbf{T}$  is minimal if no topic can be removed without losing representativeness.

# Topical Sequence Profiling

## Problem Statement (continued)

A topic embedding  $\mathbf{T}$  of  $\mathcal{D} = \{D_1, \dots, D_n\}$ :

$$\mathbf{T} = \begin{matrix} \text{Cluster analysis} \\ \vdots \\ \text{Retrieval model} \end{matrix} \left[ \begin{array}{ccc} D_1 & \dots & D_n \\ 4\% & & 40\% \\ & \dots & \\ 7\% & & \mathbf{T}_{kn} \end{array} \right]$$


---

Coverage 72% ...  $C_n$

### 3. Minimality.

A topic embedding  $\mathbf{T}$  is minimal if no topic can be removed without losing representativeness.

Sought is the minimum representative topic embedding  $\mathbf{T}^*$  with the highest mean topical variance.  $\mathbf{T}^*$  is called the topical sequence profile of  $\mathcal{D}$ .

# Topical Sequence Profiling

## Approach

Our approach comprises two steps:

### (1) *Topic Acquisition.*

- Acquisition and assignment of topics to documents.
- Result: a representative but not minimal topic embedding  $\mathbf{T}$ .
- We consider Wikipedia articles as topics with their titles as explicit labels.

Operationalization: ESA-inspired model with heuristic candidate search.

### (2) *Topic Selection.*

- Find  $\mathbf{T}^*$  given  $\mathbf{T}$ .

Operationalization: Greedy topic selection.

# Topical Sequence Profiling

## Approach: (1) Topic Acquisition

1. Represent each **Wikipedia article** (title  $\sim$  topic) as BM25 vector  $t$ .

Represent each **document** in  $\mathcal{D}$  as BM25 vector  $d$ .

Represent the **topic domain** of  $\mathcal{D}$  as centroid  $\bar{d}$  of all vectors,  $\bar{d} = \sum_{i=1}^{|\mathcal{D}|} d_i$ .

2. For each Wikipedia article:

Compute the similarity  $\rho(t, \bar{d})$  of the <Wikipedia article, topic domain> pair.

Remember the top similar Wikipedia articles as seed documents.

3. Explore (BFS) the Wikipedia link graph starting with the seed documents.

For each unassigned document in  $\mathcal{D}$ :

Compute the similarity  $\rho(t, d)$  of the <Wikipedia article, document> pair.

If  $\rho(t, d) > \rho(t, \bar{d})$  then assign topic  $t$  to document  $d$ .



# Topical Sequence Profiling

## Approach: (1) Topic Acquisition

1. Represent each **Wikipedia article** (title  $\sim$  topic) as BM25 vector  $t$ .  
Represent each **document** in  $\mathcal{D}$  as BM25 vector  $d$ .  
Represent the **topic domain** of  $\mathcal{D}$  as centroid  $\bar{d}$  of all vectors,  $\bar{d} = \sum_{i=1}^{|\mathcal{D}|} d_i$ .
2. For each **Wikipedia article**:  
    Compute the similarity  $\rho(t, \bar{d})$  of the **<Wikipedia article, topic domain>** pair.  
    Remember the top similar Wikipedia articles as seed documents.
3. Explore (BFS) the Wikipedia link graph starting with the seed documents.  
    For each unassigned document in  $\mathcal{D}$ :  
        Compute the similarity  $\rho(t, d)$  of the **<Wikipedia article, document>** pair.  
        If  $\rho(t, d) > \rho(t, \bar{d})$  then assign topic  $t$  to document  $d$ .

# Topical Sequence Profiling

## Approach: (1) Topic Acquisition

1. Represent each **Wikipedia article** (title  $\sim$  topic) as BM25 vector  $t$ .  
Represent each **document** in  $\mathcal{D}$  as BM25 vector  $d$ .  
Represent the **topic domain** of  $\mathcal{D}$  as centroid  $\bar{d}$  of all vectors,  $\bar{d} = \sum_{i=1}^{|\mathcal{D}|} d_i$ .
2. For each **Wikipedia article**:  
    Compute the similarity  $\rho(t, \bar{d})$  of the **<Wikipedia article, topic domain>** pair.  
    Remember the top similar Wikipedia articles as seed documents.
3. Explore (BFS) the Wikipedia link graph starting with the seed documents.  
    For each unassigned **document** in  $\mathcal{D}$ :  
        Compute the similarity  $\rho(t, d)$  of the **<Wikipedia article, document>** pair.  
        If  $\rho(t, d) > \rho(t, \bar{d})$  then assign topic  $t$  to document  $d$ .

# Topical Sequence Profiling

## Approach: (2) Topic Selection

Result of the topic acquisition step:

$$\mathbf{T} = \begin{array}{c} t_1 \\ \vdots \\ t_n \end{array} \left[ \begin{array}{ccc} D_1 & \cdots & D_n \\ \mathbf{T}_{11} & & \mathbf{T}_{1n} \\ & \ddots & \\ \mathbf{T}_{1m} & & \mathbf{T}_{mn} \end{array} \right]$$

---

100%       $\cdots$       100%

Greedy topic selection strategy:

1. Sort the topics of the embedding by ascending diversity  $\text{Var}(\mathbf{T}_{i:})$ .
2. For each topic  $t$ :  
If the embedding  $\mathbf{T}$  without  $t$  is representative, remove  $t$  from  $\mathbf{T}$ .

## Notes:

- ❑ The above greedy strategy returns a minimal topic embedding  $\hat{\mathbf{T}}^*$ —precisely, a topic embedding that can not be made smaller.
- ❑ However, the topic embedding  $\hat{\mathbf{T}}^*$  may not be minimum, i.e.,  $|\hat{\mathbf{T}}^*| < |\mathbf{T}^*|$ .
- ❑ Computing  $\mathbf{T}^*$  is an NP-hard problem. For its decision variant (existence of a  $\mathbf{T}$  for a  $k$ -bounded mean topic variance) the NP-completeness may be shown.
- ❑ Since the above strategy removes topics in ascending order of their diversity, the algorithm strives for maximizing the mean topic diversity.
- ❑ The heuristic is effective, if the topics are orthogonal (since the coverage percentages won't “add up” on few text collections).

# Topical Sequence Profiling

## Related Work

Existing approaches focus on (single) trend detection—not on coverage.\*

- (a) Apply LDA topic detection to  $\mathcal{D}$ , irrespective the topic distributions in the individual text collections. [Griffiths:2004, Hall:2008, Yeung:2011]
- (b) Apply LDA topic detection to each text collection  $D \in \mathcal{D}$  individually, then align the topics across the individual text collections. [Swan:2000, Wang:2005]
- (c) Model time explicitly as a parameter of the LDA topic model. [Blei:2006, Wang:2006]

\*) To showcase results, topics are cherry-picked or the hottest/coldest topics are presented. The extend to which the presented topics cover (= characterize) the sequence as a whole is not examined.

# Topical Sequence Profiling

## Case Study

$\mathcal{D}$  = SIGIR conference proceedings from 2007 to 2015.

	Year	# Papers
$D_1$	2007	198
$D_2$	2008	193
$D_3$	2009	193
$D_4$	2010	214
$D_5$	2011	232
$D_6$	2012	216
$D_7$	2013	205
$D_8$	2014	226
$D_9$	2015	193

# Topical Sequence Profiling

## Case Study

$\mathcal{D}$  = SIGIR conference proceedings from 2007 to 2015.

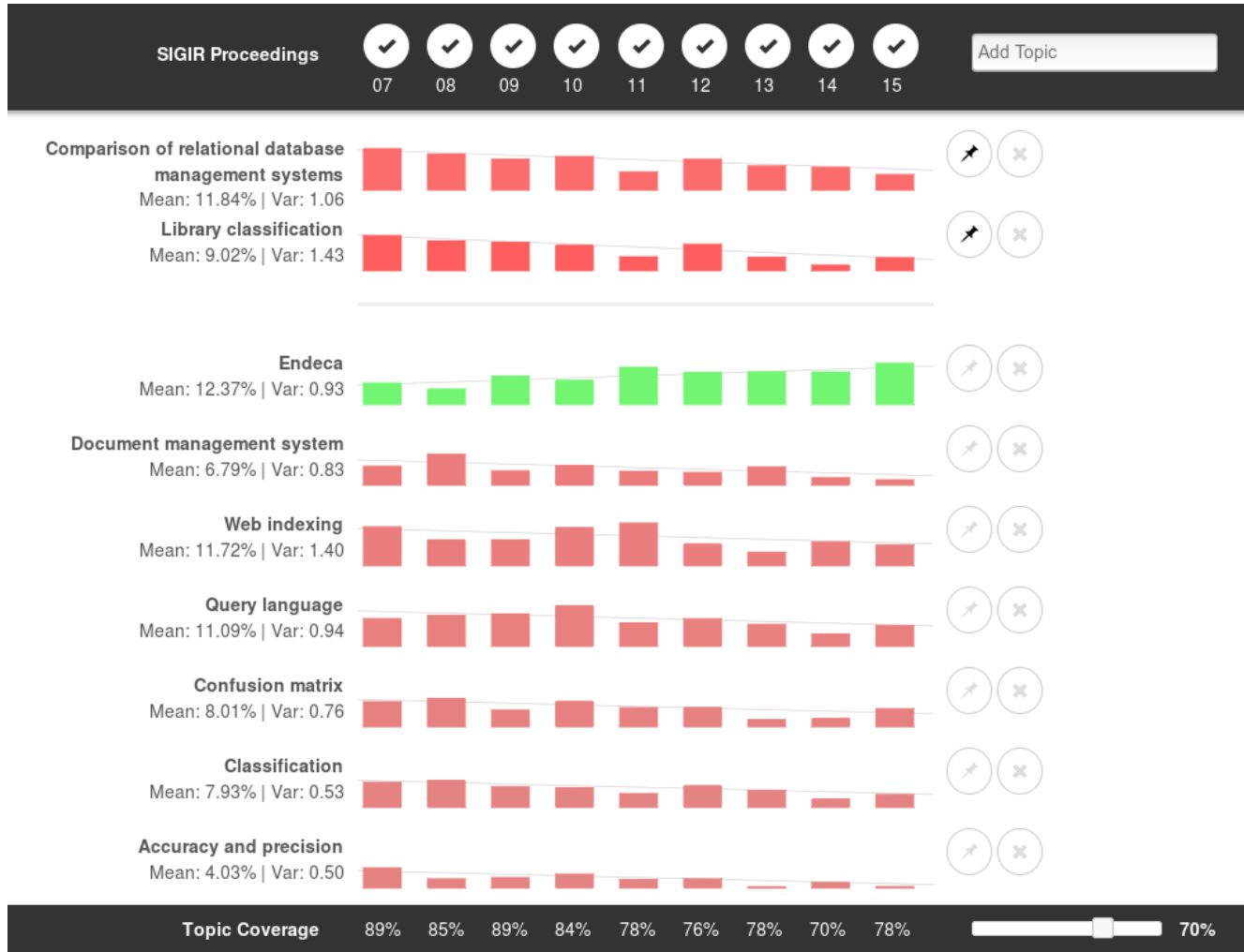
	Year	# Papers
$D_1$	2007	198
$D_2$	2008	193
$D_3$	2009	193
$D_4$	2010	214
$D_5$	2011	232
$D_6$	2012	216
$D_7$	2013	205
$D_8$	2014	226
$D_9$	2015	193

Wikipedia articles (seed documents) determined in the acquisition phase, along with the averaged similarity  $\rho(t, \bar{d})$ :

1. Concept Search (0.678)
2. Information Retrieval (0.593)
3. Human-Computer Information Retrieval (0.588)
4. Web Query Classification (0.582)
5. Enterprise Search (0.549)
6. Search engine technology (0.540)
7. Document retrieval (0.539)
8. Cognitive models of information retrieval (0.524)
9. Federated search (0.524)
10. Web search query (0.518)

# Topical Sequence Profiling

Case Study [Demo: local, web]





# Topical Sequence Profiling

## Summary and Outlook

- Topical sequence profiling means to *cover* a set  $\mathcal{D}$  of text collections.
- Such a cover should be representative, informative, and minimal; a cover is called topic embedding  $\mathbf{T}$ .
- We determine a minimal  $\hat{\mathbf{T}}^*$  within two steps:
  - (1) topic acquisition with Wikipedia, and
  - (2) heuristic topic selection by variance maximization.
- Determining the optimum cover  $\mathbf{T}^*$  is NP-hard.

What is missing + further steps:

- *Thorough* evaluation of our approach:
  - effectiveness of Wikipedia-based topic acquisition
  - comparison to LDA-based approaches
  - approximation quality of the greedy selection heuristic
- Tap the potential of user interaction

# Topical Sequence Profiling

## Summary and Outlook

- Topical sequence profiling means to *cover* a set  $\mathcal{D}$  of text collections.
- Such a cover should be representative, informative, and minimal; a cover is called topic embedding  $\mathbb{T}$ .
- We determine a minimal  $\hat{\mathbb{T}}^*$  within two steps:
  - (1) topic acquisition with Wikipedia, and
  - (2) heuristic topic selection by variance maximization.
- Determining the optimum cover  $\mathbb{T}^*$  is NP-hard.

## What is missing + further steps:

- *Thorough* evaluation of our approach:
  - effectiveness of Wikipedia-based topic acquisition
  - comparison to LDA-based approaches
  - approximation quality of the greedy selection heuristic
- Tap the potential of user interaction

Thank you for listening!

# Topical Sequence Profiling

## Related (1): Cluster Labeling

Topical sequence profiling is closely related to cluster labeling:

- (a)  $\mathcal{D}$  can be considered as a clustering of documents.  
Assign to each  $D \in \mathcal{D}$  some topic  $t$  as label.
- (b) The categorization of a set  $D \in \mathcal{D}$  can be considered as non-exclusive clustering task. Assign the topics  $t \in \mathbf{T}$  as labels to the clusters in  $D$ .

Property	View (a)	View (b)
Unique	✗	✓
Summarizing	✓	✓
Expressive	✓	✓
Discriminating	✗	✓
Contiguous	✗	✓
Irredundant	✓	✓
Representative	✓	✓
Diverse	✓	✗
Minimal	✓	✓

In topical sequence profiling the desired properties of cluster labels [Stein:2004b] depend on whether we consider view (a) or view (b).

# Topical Sequence Profiling

## Related (2): Cluster Analysis

Topic acquisition and assignment in topical sequence profiling is closely related to the principle of descriptive cluster analysis. [Hoppe:2010]