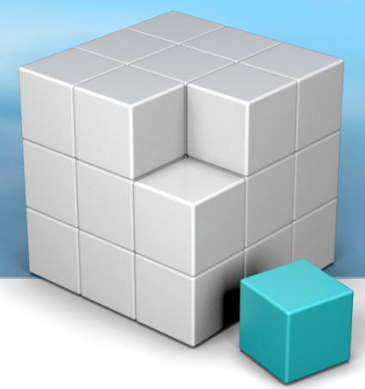# Incremental Pseudo Rectangular Organization of Information Relative to a Domain

**RAMiCS 13**
**Cambridge , UK**
**September 17 - 20**

Sahar Ismail & Ali Jaoua

1

# Agenda

1. **Motivation and Objectives**

2. **Background**

3. **Solution Model**

4. **Evaluation**

5. **Applications / Demo**

6. **Conclusion**

# Motivation

- Information is doubling roughly every 11 hours in 2011

  *Gartner and Accenture*

- Worldwide information volume is growing at a minimum rate of 59% annually                    *IBM Study*

- 70% to 85% of data is "unstructured"                    *Gartner*

- 87% of performance issues in application databases are related in some way to data growth                    *OAUG*

# ".. There Is No Such Thing as Information Overload…No Giving Up, No Surrender .."

M. Adrian

# Motivation

Solutions in market: mostly organize info by date, section
→ Manually defined tags

# Motivation

- Meta Search Engine → Too many clusters
- Not incremental, new website → do it all over again

# Motivation

- Text summarization → non incremental

# Incremental Pseudo Structuring System IPS

- **Objectives:**

  - Perform incremental information organization using pseudo maximal rectangles (new conceptual constructs)

  - Process large scale data incrementally with reasonable quality and efficiency

  - Produce domain-biased information structuring

- **Supported applications:**

  - Macro-level structuring  (Documents X Words)

  - Text Summarization (Sentences X Words)

  - Feature extraction

# Background
## **Information Representation**(1/2)

Formal Context I: (A,B,R)

| R | B1 | B2 | B3 | B4 |
|---|----|----|----|----|
| **A1** | 1 | 1 | 0 | 0 |
| **A2** | 1 | 1 | 0 | 0 |
| **A3** | 0 | 1 | 1 | 0 |
| **A4** | 0 | 1 | 1 | 1 |
| **A5** | 0 | 0 | 1 | 1 |

If **A**X**B** $\subseteq$ **A´**X**B´** $\subseteq$ **R**, and we have **A = A´** and **B = B´** ➜ the **Rectangle is maximal (Concept)**

Concepts c1 and c2 are connected iff c1 $\leq$ c2 and $\nexists$ c3 such that c1 $\leq$ c3 $\leq$ c2.

Galois Lattice

{} {A1, A2, A3, A4, A5}

{B2}{A1, A2, A3, A4}

{B3}{A3, A4, A5}

{B1, B2}{A1, A2}

{B2, B3}{A3, A4}

{B3, B4}{A4, A5}

{B2, B3, B4}{A4}

{B1, B2, B3, B4}{}

# Information Representation(2/2)

Formal Context I: (A,B,R)

| R | B1 | B2 | B3 | B4 |
|---|---|---|---|---|
| A1 | 1 | 1 | 0 | 0 |
| A2 | 1 | 1 | 0 | 0 |
| A3 | 0 | 1 | 1 | 0 |
| A4 | 0 | 1 | 1 | 1 |
| A5 | 0 | 0 | 1 | 1 |

**Minimal Coverage: {C4,C5,C6}**

**Pseudo Max. Rectangle (A4, B3)**

Is the union of all max. rectangles containing
$(A_4, B_3)$: PS = **$I(B_3.R^{-1})$ o R o $I(A_4.R)$**

Lattice of Context I

{} {A1, A2, A3, A4, A5}

{B2}{A1, A2, A3, A4}

{B3}{A3, A4, A5}

C4
{B1, B2}{A1, A2}

C5
{B2, B3}{A3, A4}

C6
{B3, B4}{A4, A5}

{B2, B3, B4}{A4}

{B1, B2, B3, B4}{}

# Solution Model
## Overview

**The solution was modeled in terms of the following dimensions:**

1. Domain of knowledge
2. Information store
3. Non incremental and incremental algorithms

Websites

Documents

News

Social Media

Unorganized sources

Uniformity in structuring

Add , Delete update

Pseudo conceptual structure

PS

PS

Ps

Ps

Ps

Ps

Ps

Structure maintained incrementally

# Solution Model
# **Domain of knowledge**

- Domain information categorized into disjoint bag of words

- Each bag represents a category and contains representative category words

- Aid labeling and structuring → more semantics but less minimality

| | |
|---|---|
| **Management Change** | • retire, appoint, resign, assign, promote, demote, hire |
| **Transaction** | • buy, sell, Lease, rent, deal, loan, contract, asset, exchange |
| **Performance** | • grow, shrink, increase, decrease, lose , gain , profit, earn, drop, jump, discount, raise |
| **Others** | • <Empty bag> |

**FWATCH Domain  Bags of words**

# Solution Model
# **Structured information store**

- Input: Binary relation representing documents
- Output of the structuring process:
  - Multi-root tree with category nodes at the first level
  - Subsequent levels under each category hold pseudo max. rectangles organized as a heap

| **Management Change** | **Transaction** | **Performance** | **Others** |
|---|---|---|---|

**Pseudo concept 1**

**Pseudo concept 3**

**Pseudo concept 4**

**Pseudo concept 5**

**Pseudo concept 2**

**Pseudo concept 7**

**Pseudo concept 6**

# Solution Model
# **Algorithms**

- ✓ Find minimal pseudo rectangular coverage non-incrementally + structure relevant to domain

- ✓ Maintain the structure incrementally in acceptable **quality** and reasonable **efficiency**

| Management Change | Transaction | Performance | Others |
|---|---|---|---|

Pseudo MAxRect1    Pseudo MaxRect3    Pseudo MaxRect4    Pseudo MaxRect5

Pseudo MaxRect2    Pseudo MaxRect7     Pseudo MaxRect6

# **Non-incremental algorithm** (1/3)

## **1.** Pre-preprocessing + NLP → Create Formal Context

**3 Documents**

.. these reflect **revenue** volatility ~~because of~~ ~~the~~ link of pro prices rates

**Simulate** results ~~to~~ interpreted in terms of ~~the~~ sks nd n

determine unit **costs** ~~and~~ **revenue** prices ~~for~~ ~~the~~ projects. Unit costs ~~are~~ presented ~~for~~ the projects ~~on~~ apriority basis, ~~without~~ taking external ~~and~~ actual project arrangements ~~into~~ account. ~~As~~ provided ~~in the~~ **simulation**

| Pair | weight |
|------|--------|
| **(1,1)** | 2 |
| **(2,2)** | 3 |
| **(3,1)** | 5 |
| **(3,2)** | 1 |
| **...** | ... |
| **(3,n)** | -1 |

**R[3][n]: 3 documents associated to n non-empty words**

| R | revenue | simulate | ... | cost |
|------|---------|----------|-----|------|
| **doc1** | 1 | 0 | ... | 0 |
| **doc2** | 0 | 1 | ... | 0 |
| **doc3** | 1 | 1 | .. | 1 |

# Methodology
# Non-incremental algorithm (2/3)

## 2. Calculate minimal pseudo-conceptual coverage

**2.1** Sort pairs in descending order of their weights

| Pair | weight |
|------|--------|
| ~~(3,1)~~ | ~~5~~ |
| ~~(2,2)~~ | ~~3~~ |
| ~~(1,1)~~ | ~~2~~ |
| ~~(3,2)~~ | ~~1~~ |
| ... | ... |
| ~~(3,n)~~ | ~~1~~ |

**2.2** Highest weight pair ( Familiar with respect to domain)

**2.6** Next uncovered highest weight pair

**2.3** find
PS (3,1) =I(1.$R^{-1}$) o R o I(3.R)

= {(3,1),(3,2), (3,n) , (1,1)}

← **2.4** Mark all pairs as covered

← **2.8** Mark all pairs as covered

**2.7** PS (1,1) =I(1.$R^{-1}$) o R o I(1.R)

= {(1,3),(1,1)}

**2.5** Select category and label

**Category is the most overlapping with ps words → label is the best not selected**

| Category 1 | Category 2 |
|------------|------------|
| PS (3,1) | PS (2,2) |

Stop when all pairs are covered

# Methodology
## Non-incremental algorithm (3/3)

- **High degree of scalability**
  - Number of pseudo-concepts is small compared to size of concepts in a lattice or a minimal coverage set (exponential)
  - Number of pseudo concepts is bound by the number of pairs in the relation
- **Efficient**
  - Worst case time complexity is in **O ($N^2$)** ; N is density of 1's in R
  - Best case time complexity is in **O(N log N)**

- Resulting structure is **domain sensitive**

# Solution Model
## Incremental Structuring Algorithm

**The new methods cover the following update cases:**

- Addition and deletion of domain elements
- Addition and deletion of codomain elements
- Addition and deletion of associations
- Updates to the domain of knowledge

# Solution Model
# **Incremental Add Algorithm** (1/4)

**Algorithmic skeleton**

1. Identify parts of R to update on the addition of a pair or a domain or codomain element

2. Identify pseudo maximal rectangles to be updated with the new information

3. Structure information identified as irrelevant to all existing pseudo maximal rectangles in a similar fashion to the non-incremental algorithm

# **Incremental Add Algorithm** (2/4)

1. **What to update in R**

On the arrival of new information , some pairs will require weight recalculation

$$s(\mathbf{d}, w) = \left(|\mathbf{d}.R| \times |w.R^{-1}|\right) - \left(|\mathbf{d}.R| + |w.R^{-1}|\right)$$

→ Overhead increases as the cardinality of the update set increase

New domain element

|   | a | b | c |
|---|---|---|---|
| 1 | 1 | 1 | 0 |
| 2 | 1 | 1 | 0 |
| 3 | 1 | 0 | 0 |
| 4 | 0 | 1 | 1 |
| 5 | 1 | 0 | 1 |
| x | 1 | 1 | 0 |

New codomain element

|   | a | b | c | y |
|---|---|---|---|---|
| 1 | 1 | 1 | 0 | 1 |
| 2 | 1 | 1 | 0 | 0 |
| 3 | 1 | 0 | 0 | 0 |
| 4 | 0 | 1 | 1 | 1 |
| 5 | 1 | 0 | 1 | 0 |

New Pair

|   | a | b | c |
|---|---|---|---|
| 1 | 1 | 1 | 0 |
| 2 | 1 | 1 | 0 |
| 3 | 1 | 0 | 0 |
| 4 | 0 | **1** | 1 |
| 5 | 1 | 0 | 1 |

# Incremental Add Algorithm (3/4)

## 2. What to update in PS tree

Relevant pairs → Find home from existing pseudo max. rectangles

Not relevant pairs → Structure as new

# Methodology
## **Incremental Add Algorithm** (4/4)

- Different addition scenarios:
  - **O (n log n)** in the best case and **O (n2 + N log N)** in the worst case ;        (n : count of new pairs, N: count of all pairs in R)

- Add algorithms perform better when the Context is stable or  n << N ;                              (N is count of all pairs in R)

- Weight recalculation ➔ Big computational overhead
  - Relax the immediate weight update condition
  - Drawback: Info store isn't up to date at all times (update strategy)

# Evaluation
## **Approach and metrics**(1/3)

- Incremental vs. non-incremental methods
- Metrics:
  - Cohesion: measures inner tightness in a pseudo concept
  - Separation: measures inter-pseudo concept overlapping
  - and performance (incremental vs. non incremental)

- Testing data: ~1020 documents from the NSF Research Award Abstracts 1990-2003

- Bag-of-words used represent the entire vocabulary of the documents in use

# Evaluation
## **Approach and metrics**(2/3)

- Quality of structure – **Cohesion**
  - Measures the degree of tightness within a pseudo-concept (how related the contained pairs are)

$$\text{Cohesion} = \sum_{i=1}^{k} \text{Cohesion}(\text{PSi})$$

$$\text{Cohesion}(\text{PSi}) = \frac{|\text{PSi}|^2}{|\text{Ai}| \times |\text{Bi}|}$$

- High cohesion is favorable , higher inner tightness means
  - Less pseudo maximal rectangles to maintain
  - Better scalability
  - Perfect cohesion ➜ Pseudo maximal rectangle is a maximal rectangle

# Evaluation
## Approach and metrics(3/3)

- Quality of structure – **Separation**
  - Measures how independent or distinct two pseudo maximal rectangles are

$$\text{Separation}\,(PSi, PSj) = \frac{1}{|A_i \cap Aj| \times |B_i \cap Bj|}$$

$$\text{Separation} = \sum_{j=0}^{k-1} \sum_{j=i+1}^{k} \text{Separation}\,(PSi, PSj)$$

- Separation is favorable for incremental management, less overlapping means:
  - less update requirements in case of change
  - more stable structure

# Evaluation – Incremental Add
## Experimentation results - Cohesion



- Cohesion is increasing as more documents are added incrementally
- Incremental-add algorithm favors updating existing relevant pseudo concepts over creating new pseudo concepts
- Relevant pairs are added which increases the degree of completeness of pseudo concept and increase its size

# Evaluation – Incremental Add
## Experimentation results - Separation



- Separation is declining compared to the non-incremental method as more documents are added to the context

- New documents added update all relevant pseudo concepts which decreases separation → workaround: Don't update all

# Evaluation – Incremental Add
## Experimentation results - Performance



**Legend:** None incremental Time · Incremental time

*Y-axis:* Run time (ms) — 0 to 140000
*X-axis:* Size of the Context (Density of 1's) — 62000 to 63000

- Time difference between the two approaches is considerable

- Both grow as a function of the density of 1's in the context

- Counting the accumulative time for obtaining a structure using the incremental method shows it isn't suited for structuring a space from scratch

# Conclusion

- IPS System uses new methods to handle incremental changes in an information store

- New conceptual constructs were used → Pseudo maximal rectangles

- Minimal pseudo coverage provides high level summary → promising scalability.

- Domain information is used to reduce the linguistic noise and improve labeling process

- System can handle incremental corpus organization, text summarization and feature-extraction

# Applications
## Macro-level structuring (Documents X Words)

## Macro structuring and Feature Extraction

**Pseudo Coverage** (folder: D:\test\web)

open all | close all

- Structure
  - Management Change
    - director
      - appoint
        - manag
        - grow
      - execut
        - nate
        - profit
  - Transaction
  - Performance
    - rose
  - Others

**Show pairs** ✔

**appoint**

appoint , Automot , Group ,
appoint ,
CEO , appoint , chairman , CEO , point , appoint , senior , senior , point , manag , board , grow , Ernst , compani , Inc ,
Chief , Execut , Officer , ad , Presid , global , rise , AIDS , Found , Pendleton ,
CEO , contin , expand , rise , Comium ,
Ramada , Muscat ,
grow , increas ,
ad ,
appoint , board , ad , expand ,
Chief , Execut , CEO , execut , appoint , Al ,
manag , add , expand ,
grow , edg ,
senior , appoint , Etihad ,
Officer , Chief ,
appoint , Dubai ,
appoint ,
Chief , Execut , CEO , grow ,
appoint ,
resign ,
Execut , Chief , Non ,

**Show words** ✔

period bank effect connect compon busi billion appoint execut director current vice presid parti secretari gener automot
group hold compani limit control sharehold brillianc china entitl annual remuner announc wholli own subsidiari automobil

# Applications
## Text Summarization (Sentences X Words)

**Micro Pseudo Coverage**

(File: D:\test\web\mana-12.txt)

**Parent pseudo concept**
- Micro Structure
  - Management Change
    - IAVI
      - AIDS
        - expect
        - Officer
      - org
        - profit
  - Transaction
  - Performance
    - rise
  - Others

open all | close all

The International AIDS Vaccine Initiative (IAVI) congratulates its founder, President and Chief Executive Officer Seth Berkley on his appointment as Chief Executive Officer of the Global Alliance for Vaccines and Immunization (GAVI). 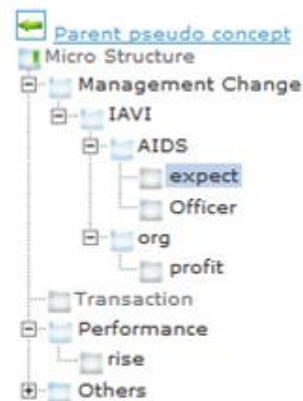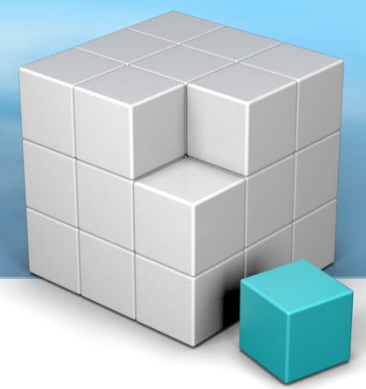The Board of GAVI, an organization that increases access to immunization in poor countries, announced the appointment today."This is a terrific opportunity for Seth and a logical next step, and we wish him well," said IAVI Board Chairman Paul Klingenstein.Berkley founded IAVI in 1996.Since then, the organization has mobilized US$ 897 million in resources for AIDS vaccine R&D, advocated to keep the development of an AIDS vaccine on the international agenda, established and maintained a network connecting clinical research centers in Africa and India, tested nine AIDS vaccine candidates in clinical trials in 11 countries and made significant contributions toward solving the problem of how to design vaccines that elicit antibodies that neutralize a broad range of variants of HIV."It was Seth's vision and energy that enabled IAVI to prosper and develop into the mature organization that it is, and he will leave it in excellent condition," said Klingenstein."It's no longer a question of whether there will be an AIDS vaccine but when, and IAVI is in an excellent position to help advance the science to the goal of an effective vaccine accessible to all."I'm excited about the new opportunity at GAVI, but of course extremely saddened to be leaving an organization so close to my heart.It's bittersweet to leave IAVI when AIDS vaccine science is showing more promise than ever.I look forward to the time when HIV will be added to GAVI's list of available vaccines," said Berkley."It's been a tremendous honor to have worked alongside all the outstanding individuals who have taken up the mission of an AIDS vaccine.I especially want to thank IAVI's generous donors, the trial volunteers and our scientific collaborators."Seth has been an indefatigable and enormously effective champion for AIDS vaccine development," said David Cook, Chief Operating Officer at IAVI."We expect he will continue to be a strong advocate from his new position at GAVI.Berkley is expected to remain in his position as President and CEO of IAVI through June.The IAVI Board of Directors will define a transition plan."Given the momentum in AIDS vaccine development and IAVI's contributions to the current excitement, a transition presents the opportunity to invigorate the organization with another outstanding leader," said Klingenstein."The Board looks forward to rising to that challenge. About IAVI The International AIDS Vaccine Initiative (IAVI) is a global not-for-profit organization whose mission is to ensure the development of safe, effective, accessible, preventive HIV vaccines for use throughout the world.Founded in 1996 and operational in 25 countries, IAVI and its network of collaborators research and develop vaccine candidates.IAVI's financial and in-kind supporters include the Bill & Melinda Gates Foundation, the Foundation for the

# Future work

- Use Ontology aligned structuring techniques

- Use higher abstractive approximation of context coverage (Union of maximal rectangles) for obtaining a higher semantic level in the structure at a relatively low cost → Hyper rectangles

- Define optimization factors for identifying optimal update mechanisms(incremental vs. non-incremental)

# Thank You
# Q & A

Incremental Pseudo Rectangular Organization of Information Relative to a Domain