QATAR UNIVERSITY

# Hyper-Rectangular Relation Decomposition

# Application for Text Mining and Structuring

**Prepared by: Fethi FERJANI**

# **Agenda**

1. **Problem and approach**

2. **About Hyper Rectangles**

3. **Hyper Rectangular Decomposition and Coverage**

4. **Conclusion and future work**

# Problem and Approach

- **What**
  - Achieve structuring and browsing capabilities in a large-scale corpus on two different levels:
    - Macro-level browsing: deals with the global visual and logical structure of the corpus (e.g. categories, sub-categories, sections, fields, etc).
    - Micro-level browsing: is used for navigating through a textual document. For example, company information, location, etc) .
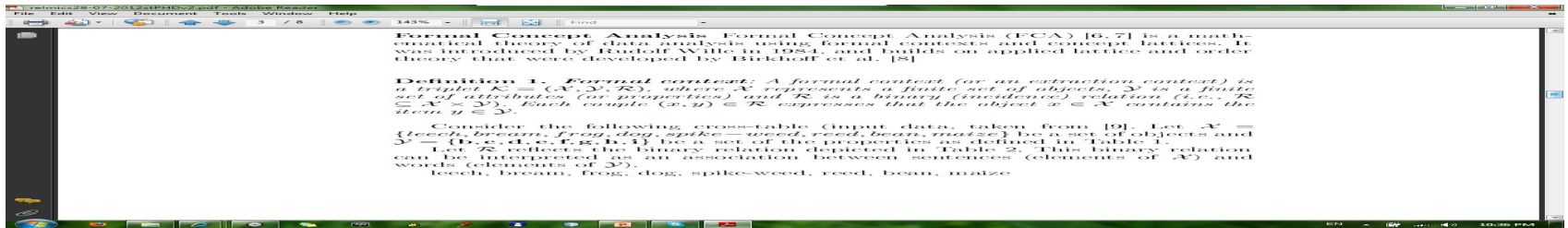
- **How**
  - New method handles construction of a Hyper Rectangular Generic tree (GT).
  - Introducing scalability and strong sematic structures.
  - Introduces a solution to the rectangles labeling challenge

# Hyper Rectangle Definition

## Formal Context



## Hyper Rectangle

**Definition 2.** (HYPER RECTANGLE)

Let $(\mathcal{X}, \mathcal{Y}, \mathcal{R})$ be a formal context and $a \in \mathcal{Y}$, an arbitrary attribute. The Hyper Rectangle, denoted by $\mathcal{H}_a$, is a sub-relation of $\mathcal{R}$ such that $\mathcal{H}_a(\mathcal{R}) = \mathcal{I}(a.\mathcal{R}^{-1}) \circ \mathcal{R}$.

ø Hyper Rectangle is the union of all possible rectangles in a binary relation involving an element of the domain or the range of a given binary relation

# Example of Hyper Rectangle

b: lives in water
d: needs chlorophyll to produce food,
f: one seed leaf,
h: has limbs,
c: lives on land,
e: two seed leaves,
g: can move around,
i: suckles its offspring.

**Table 1.** Properties Description

| | b | g | h | c | i | d | f | e |
|---|---|---|---|---|---|---|---|---|

b: lives in water
d: needs chlorophyll to produce food,
f: one seed leaf,
h: has limbs,
c: lives on land,
e: two seed leaves,
g: can move around,
i: suckles its offspring.

**Table 1.** Properties Description

**Binary Relation R**

| | b | g | h | c | i | d | f | e |
|---|---|---|---|---|---|---|---|---|
| $s_0 = leech$ | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| $s_1 = bream$ | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 |
| $s_2 = frog$ | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 |
| $s_3 = dog$ | 0 | 1 | 1 | 1 | 1 | 0 | 0 | 0 |
| $s_4 = spike - weed$ | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 1 |
| $s_5 = reed$ | 1 | 0 | 0 | 1 | 0 | 1 | 1 | 0 |
| $s_6 = bean$ | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 |
| $s_7 = maize$ | 0 | 0 | 0 | 1 | 0 | 1 | 1 | 0 |

**Table 2.** Binary Relation $\mathcal{R}$

Before providing a framework of the new definitions of HYPER RECT-ANGLE associated to a given attribute and its REMAINING RELATION, it is important to provide an example. If we consider that the attribute b has a maximum weight, we split the working binary relation in two sub-relations: The first contains all rows validating the attribute b and a second sub-relations where the attribute b is not valid. Hence, we present the following example:

**Hyper Rectangle Hb(R)**

| | b | g | h | c | i | d | f | e | | b | g | h | c | i | d | f | e |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $s_0$ | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | | | | | | | | | |
| $s_1$ | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | $s_3$ | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 1 |
| $s_2$ | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | $s_6$ | 0 | 0 | 1 | 1 | 0 | 1 | 1 | 0 |
| $s_4$ | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | $s_7$ | 0 | 0 | 1 | 1 | 0 | 1 | 0 | 0 |
| $s_5$ | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | | | | | | | | | |

**Table 3.** (**Right**): The Hyper Rectangle associated to the attribute b (**Left**): The Remaining Binary relation associated to the attribute b

# Hyper Rectangle Weighting

ø  Optimization criterion for the decomposition process

ø  Associates a strength measure to each property

ø  Higher weight indicates generality of an attribute as it is shared with more objects

ø  **Hyper Rectangle Weight**

For each HYPER RECYTANGLE, we associate a weight which measures its strengthens in terms of associations between objects and properties. The maximal weight attribute is the more general one since it is shared (directly or indirectly) by the majority of objects. A formalization of this weight is given as below.

**Definition 3.** (HYPER RECTANGLE WEIGHT)
Let $H_a(R)$ a Hyper Rectangle associate to an attribute $a$. The weight $W(H_a(R))$ of the Hyper Rectangle $H_a(R)$ is defined by a significant information optimization criteria used in the former work of [10]. The economy of a binary relation or the gain in storage space can be calculated in accordance to the following equation:

$$W(H_a(R)) = (r/(d * c)) * (r - (d + c)) \qquad (1)$$

where $r$ is the cardinality of $H_a(R)$ (i.e. the number of pairs in binary relation $H_a(R)$, $d$ is the cardinality of $Dom(H_a(R))$, and $c$ is the cardinality of $Cod(H_a(R))$.

# Maximal Hyper Rectangles

- **Maximal Hyper Rectangle**



- **Remaining Binary Relation**

# Maximal Hyper Rectangle and Remaining Relation

|   | b | g | h | c | i | d | f | e |
|---|---|---|---|---|---|---|---|---|

b: lives in water
d: needs chlorophyll to produce food,
f: one seed leaf,
h: has limbs,

c: lives on land,
e: two seed leaves,
g: can move around,
i: suckles its offspring.

**Table 1.** Properties Description

|   | b | g | h | c | i | d | f | e |
|---|---|---|---|---|---|---|---|---|
| $s_0 = leech$ | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| $s_1 = bream$ | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 |
| $s_2 = frog$ | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 |
| $s_3 = dog$ | 0 | 1 | 1 | 1 | 1 | 0 | 0 | 0 |
| $s_4 = spike - weed$ | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 1 |
| $s_5 = reed$ | 1 | 0 | 0 | 1 | 0 | 1 | 1 | 0 |
| $s_6 = bean$ | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 |
| $s_7 = maize$ | 0 | 0 | 0 | 1 | 0 | 1 | 1 | 0 |

**Table 2.** Binary Relation $\mathcal{R}$

Before providing a framework of the new definitions of HYPER RECT-ANGLE associated to a given attribute and its REMAINING RELATION, it is important to provide an example. If we consider that the attribute b has a maximum weight, we split the working binary relation in two sub-relations: The first contains all rows validating the attribute b and a second sub-relations where the attribute b is not valid. Hence, we present the following example:

|   | b | g | h | c | i | d | f | e |   |   | b | g | h | c | i | d | f | e |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $s_0$ | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |   |   |   |   |   |   |   |   |   |   |
| $s_1$ | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 |   | $s_3$ | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 1 |
| $s_2$ | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 |   | $s_6$ | 0 | 0 | 1 | 1 | 0 | 1 | 1 | 0 |
| $s_4$ | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |   | $s_7$ | 0 | 0 | 1 | 1 | 0 | 1 | 0 | 0 |
| $s_5$ | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 0 |   |   |   |   |   |   |   |   |   |   |

**Table 3.** (**Right**): The Hyper Rectangle associated to the attribute b (**Left**): The Remaining Binary relation associated to the attribute b

|   | b | g | h | c | i | d | f | e |
|---|---|---|---|---|---|---|---|---|
| $s_0$ | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| $s_1$ | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 |
| $s_2$ | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 |
| $s_4$ | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| $s_5$ | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 0 |

**Table 3.** (**Right**): The Hyper Rectangle associated to the attribute b (**Left**): The Remaining Binary relation associated to the attribute b

# Hyper Rectangular Coverage

ø Calculate the weight for all attributes and sort them descending

ø Perform generalization of the best attribute into Hyper Rectangle

ø Find the Hyper rectangles in the remaining relation recursively until the entire relation is covered  or Rm(R) is a Hyper Rectangle

ø For a given level **i**, we compute the attributes for the level **i+1** with a conservation of the coverage

▸ 9of the previous level.
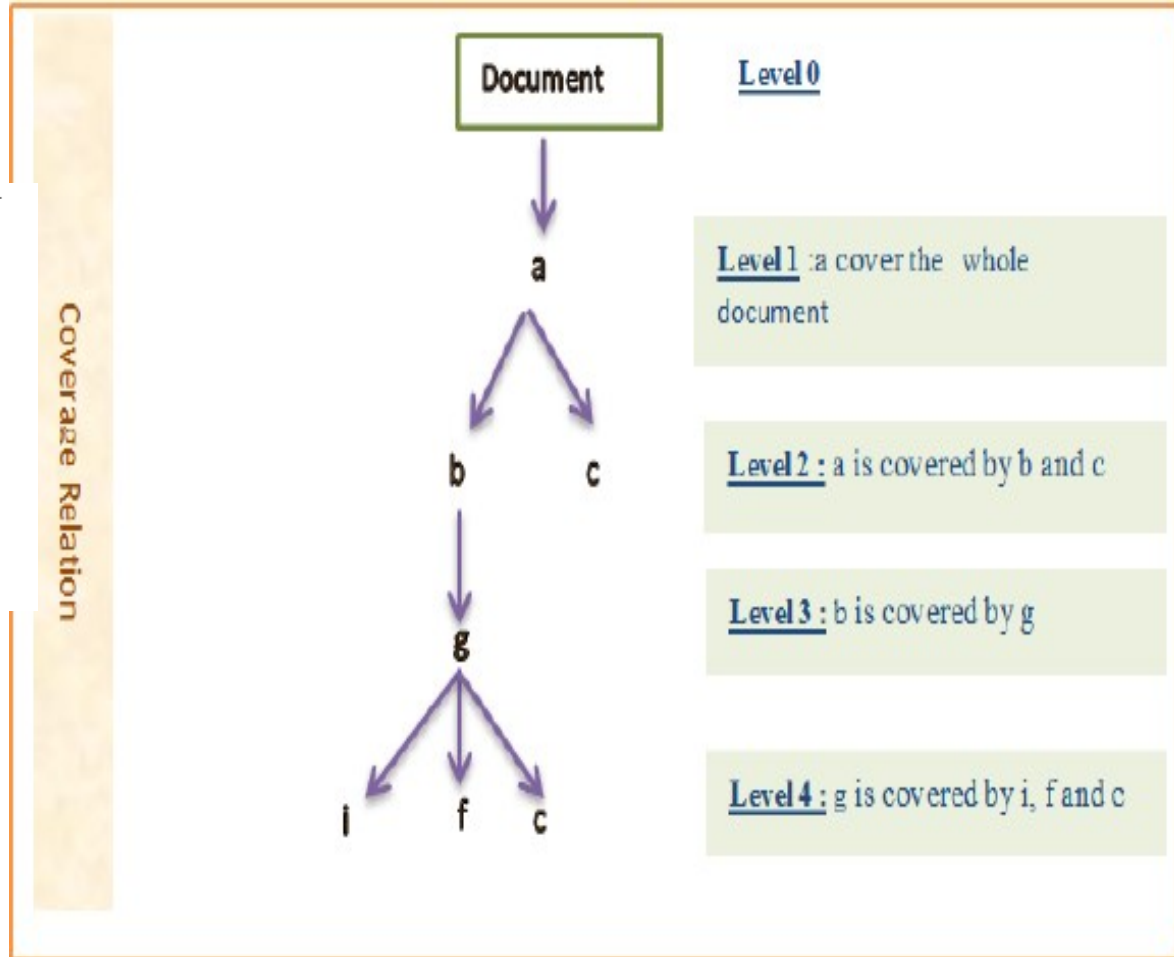
# Generic Tree of words

Before providing a framework of the new definitions of HYPER RECTANGLE associated to a given attribute and its REMAINING RELATION, it is important to provide an example. If we consider that the attribute b has a maximum weight, we split the working binary relation in two sub-relations: The first contains all rows validating the attribute b and a second sub-relations where the attribute b is not valid. Hence, we present the following example:

|     | b | g | h | c | i | d | f | e |     | b | g | h | c | i | d | f | e |
|-----|---|---|---|---|---|---|---|---|-----|---|---|---|---|---|---|---|---|
| $s_0$ | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |     |   |   |   |   |   |   |   |   |
| $s_1$ | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | $s_3$ | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 1 |
| $s_2$ | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | $s_6$ | 0 | 0 | 1 | 1 | 0 | 1 | 1 | 0 |
| $s_4$ | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | $s_7$ | 0 | 0 | 1 | 1 | 0 | 1 | 0 | 0 |
| $s_5$ | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 0 |     |   |   |   |   |   |   |   |   |

**Table 3.** (**Right**): The Hyper Rectangle associated to the attribute b (**Left**): The Remaining Binary relation associated to the attribute b
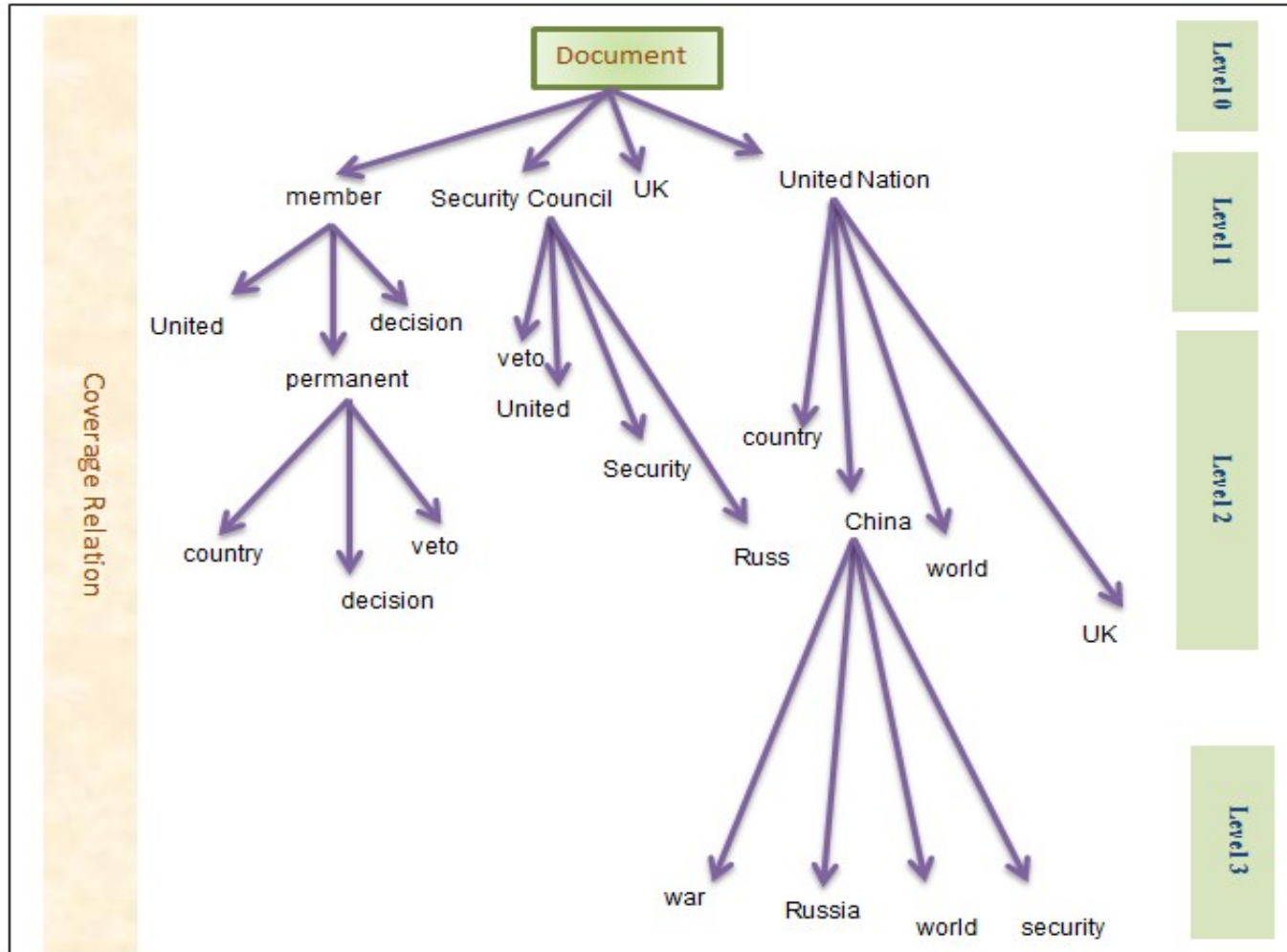
Document

Level 0

a

Level 1 : a cover the whole document

b          c

Level 2 : a is covered by b and c

g

Level 3 : b is covered by g

i     f     c

Level 4 : g is covered by i, f and c

Coverage Relation

# Illustrative Example  United Text document

The United Nations was formed at the end of the Second World War on October. It The United Nations was formed to try to make sure that future generations would not suffer from horrible damage caused by wars. It The United Nations was also intended to develop friendly relations between countries around the world and to help with any problems countries might have. The most important part of the United Nations has always been the Security Council. The Security Council was formed to try to keep peace and security around the world. In the beginning

the Security Council consisted of eleven countries; now it the Security Council has fifteen. Five of these are permanent members. Today they are Russia China France the USA and the UK. The other ten member

countries are elected for a period of two years. When decisions have to be made each member of the council has one vote.

The decision is only agreed if at least nine of the fifteen members vote for it decision. Even then any one of the five permanent members can block or stop the decision by using its veto. Only the permanent members

have a veto. This makes them very powerful.

**Table 6.** A snapshot of a United Nations textual document.

# Generic Tree Correspondent

# Conclusion

ø New heuristic based approach utilizing Hyper Rectangles for navigating and browsing inside a corpus was presented

ø Hyper Rectangle is the union of all possible rectangles in a binary relation involving an element of the domain or the range of a given binary relation

ø Hyper Rectangles are very appropriate for large data sets and may overcome the previous approaches drawbacks regarding scalability

ø As a perspective, we envisage to explore the incremental version of this approach and by the inclusion of the similarity between words in the textual document.

# Thank You