



Text Warehousing

2017. 5

Introduction

▶ Big Text Data

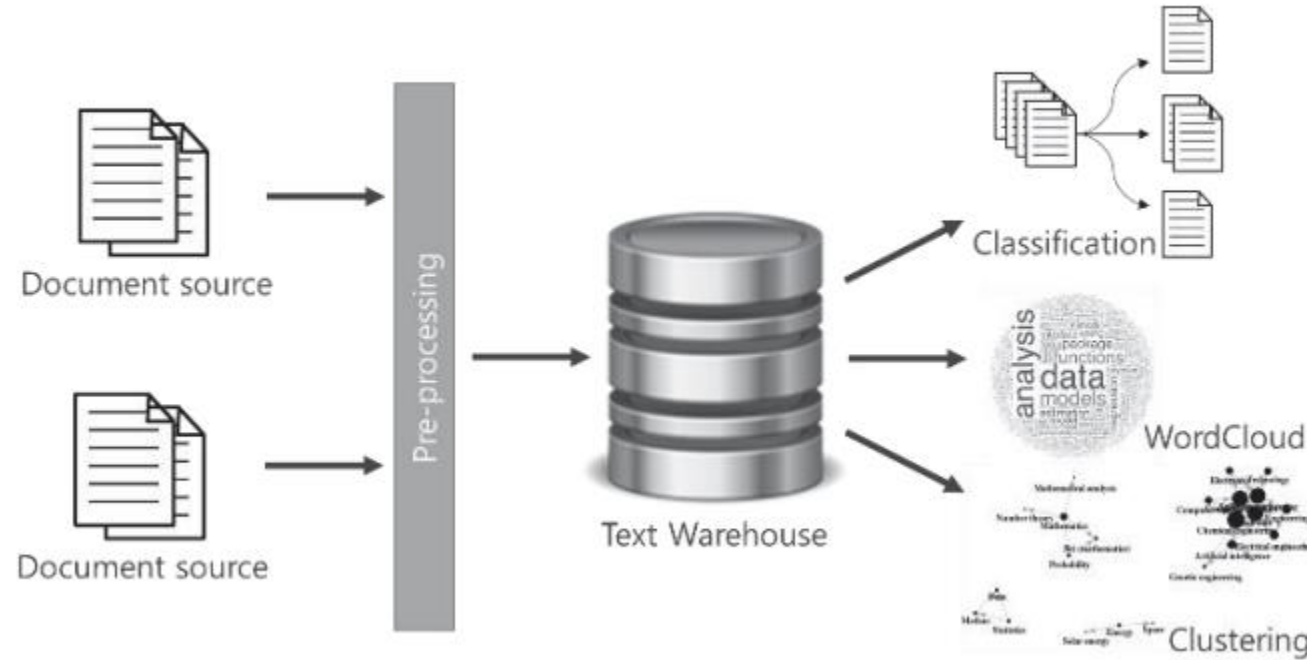
- ▶ The textual data integration in OLAP analysis and aggregating them to improve the decision-making is a challenge in Business Intelligence systems.
- ▶ important to define new aggregation techniques appropriate to textual data.

▶ Proposed technique

- ▶ IR techniques together with OLAP to better analyze textual data and extract their semantics.



Architecture



Text Warehouse 예

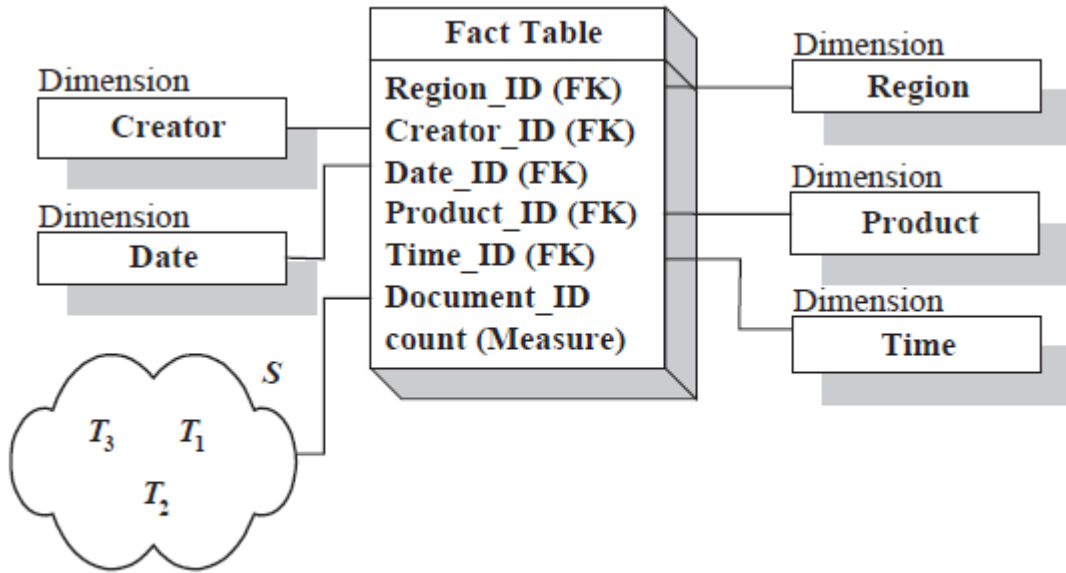


Fig. 9. An example star schema for complaint e-mail management.

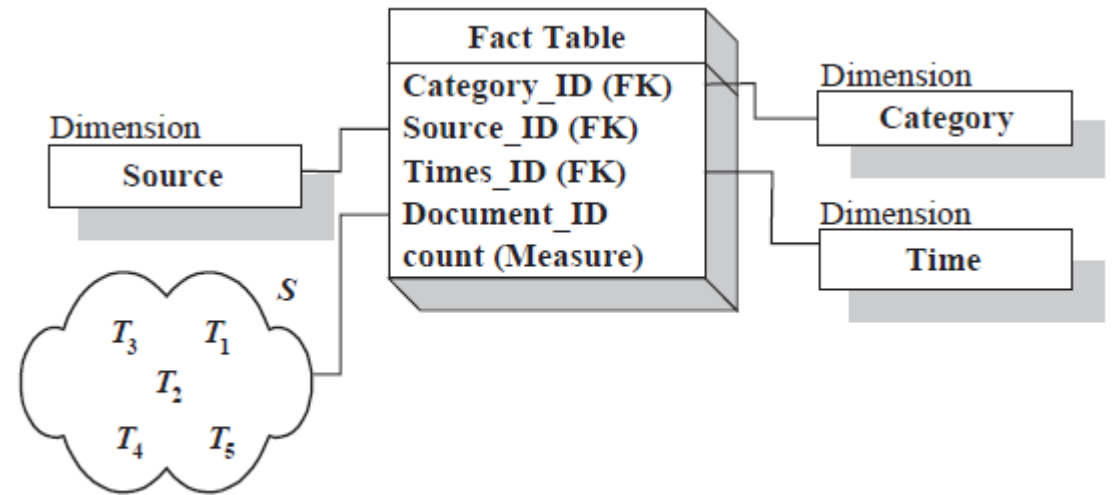
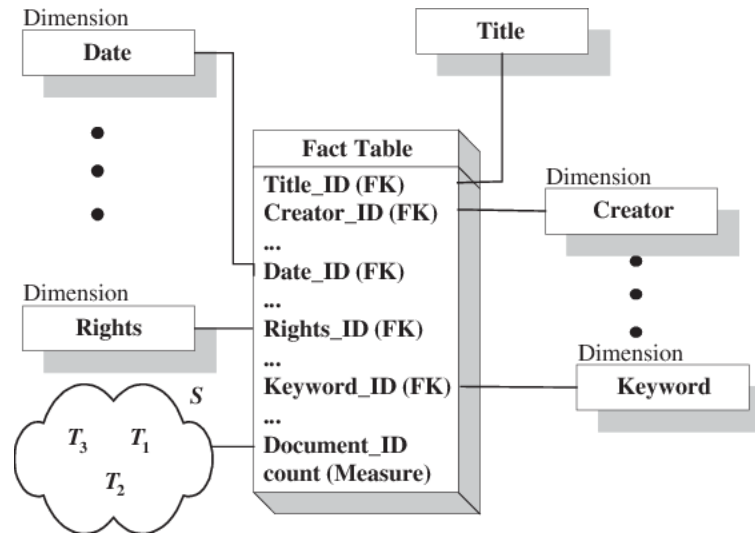


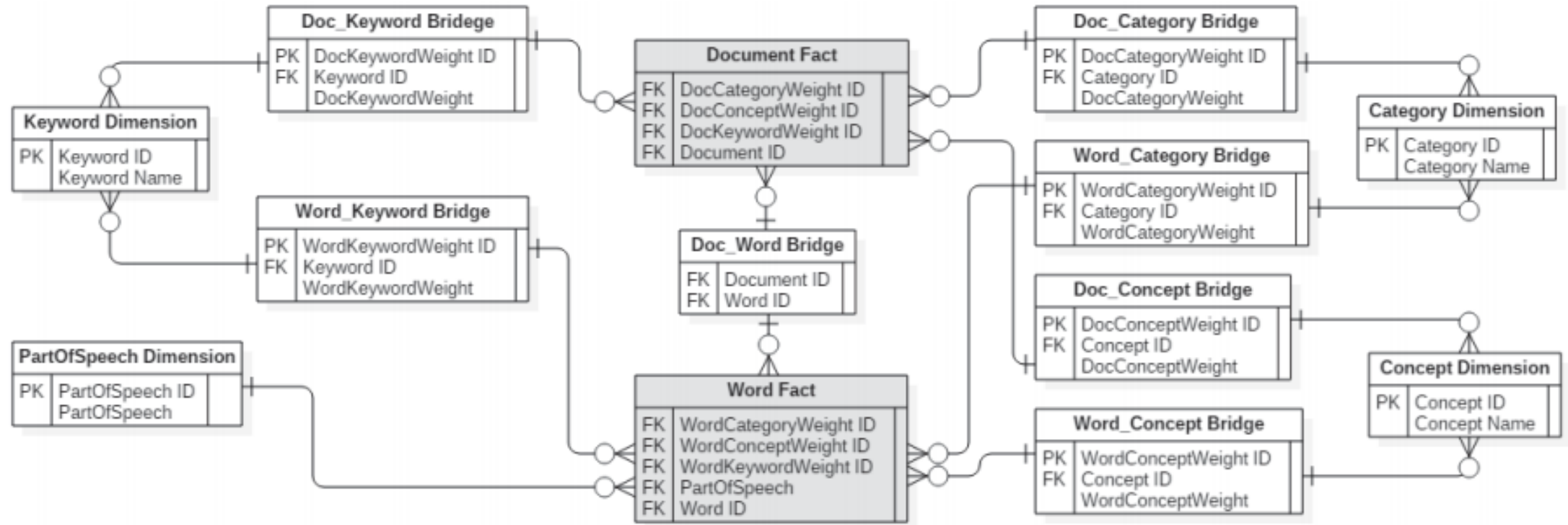
Fig. 12. An example star schema for journal paper warehousing.

Grain: document or term ?



Text Warehouse 예





Context definition in Text Warehousing

- ▶ Document context:
 - ▶ contextual factors of the analyzed documents
 - ▶ topics, domains, metadata such as title, author, date, and etc
- ▶ User context:
 - ▶ user profile which is based on his query log containing the history of queries
 - ▶ user preferences for the analysis dimensions.



Issues



- ▶ Schema design
 - ▶ Grain
 - ▶ Textual dimensions
 - ▶ Semantic dimensions
 - ▶ Bridge tables
- ▶ Text understanding
 - ▶ Concept extraction for documents and terms
 - ▶ Concept hierarchy
 - ▶ Weighting for grains



ConteXtual Text cube model: CXT-Cube

- ▶ Document representation

- ▶ Term vector

- ▶ $d = \langle w_{t_1}, w_{t_2}, w_{t_3}, \dots, w_{t_1} \rangle$

- ▶ Measure definition

- ▶ Represents each document d by several vectors of weighted concepts

$$M = \langle \overrightarrow{d_{Dim_1}}, \overrightarrow{d_{Dim_2}}, \dots, \overrightarrow{d_{Dim_*}} \rangle$$

$$\overrightarrow{d_{Dim_r}} = \langle w_{c_1}, w_{c_2}, \dots, w_{c_n} \rangle$$

Text Measures

▶ Measure definition

$$\vec{d}_{Dim_r} = \langle w_{c_1}, w_{c_2}, \dots, w_{c_n} \rangle$$

- ▶ the vector of weighted concepts of a document d in a vector space dimension specific to Dim_r and w_{c_i} is the weight assigned to the concept c_i .

▶ Example

$$M = \langle \vec{d}_{dim_z}, \vec{d}_{Dim_l}, \vec{d}_{Dim_t} \rangle$$

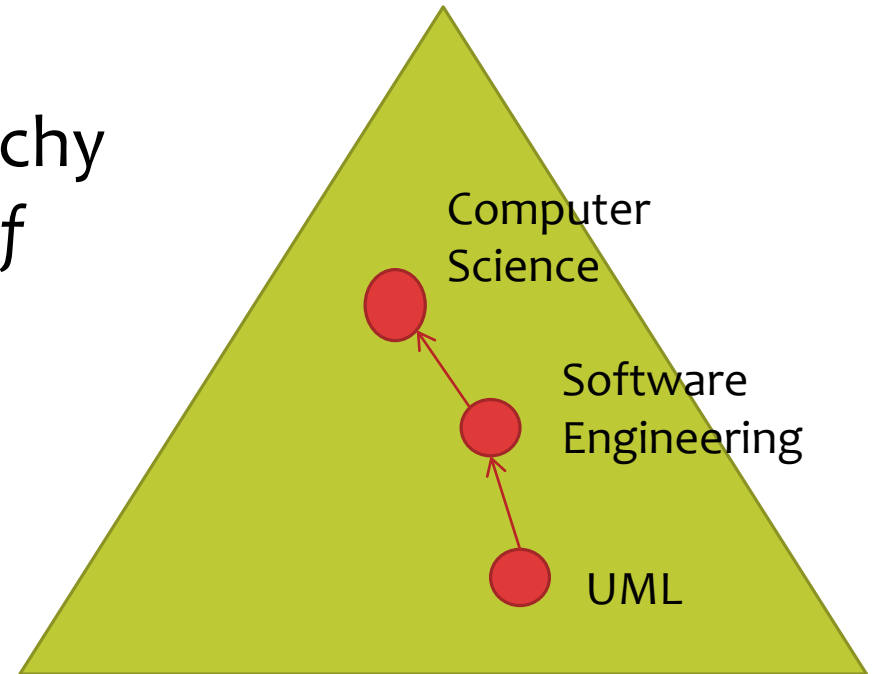
$\vec{d}_{dim_z}, \vec{d}_{Dim_l}, \vec{d}_{Dim_t}$: vectors of dimensions: TOPIC (Dim_z), LOCATION (Dim_l), and TIME (Dim_t)

Relevance Propagation

▶ Step 1

- ▶ compute the weights of the document terms which exist in the concept hierarchy by using the method Term Frequency Tf computed by the following formula:

$$Tf_{t,d} = \frac{n_{t,d}}{N_d}$$

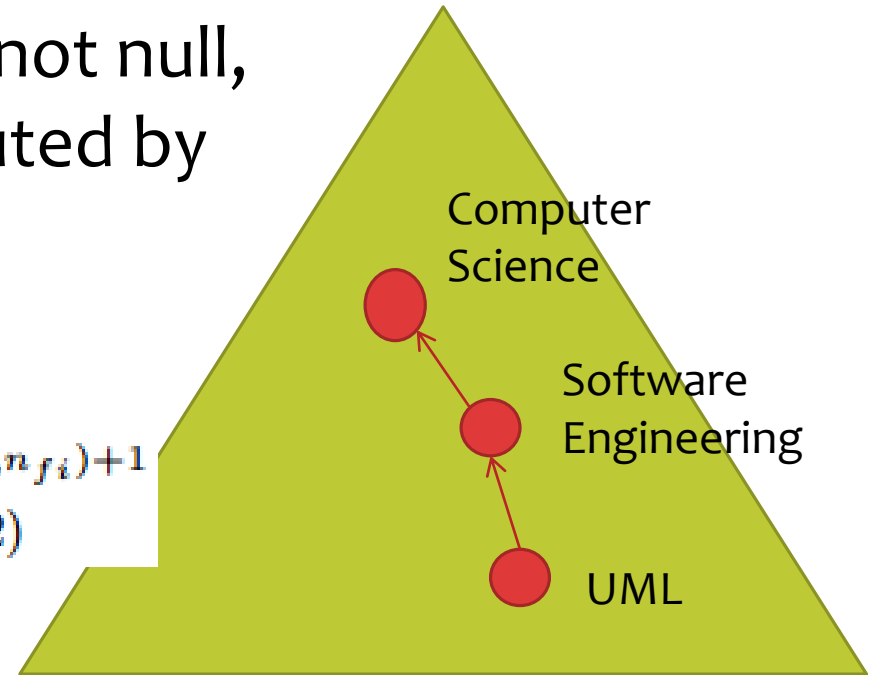


Relevance Propagation

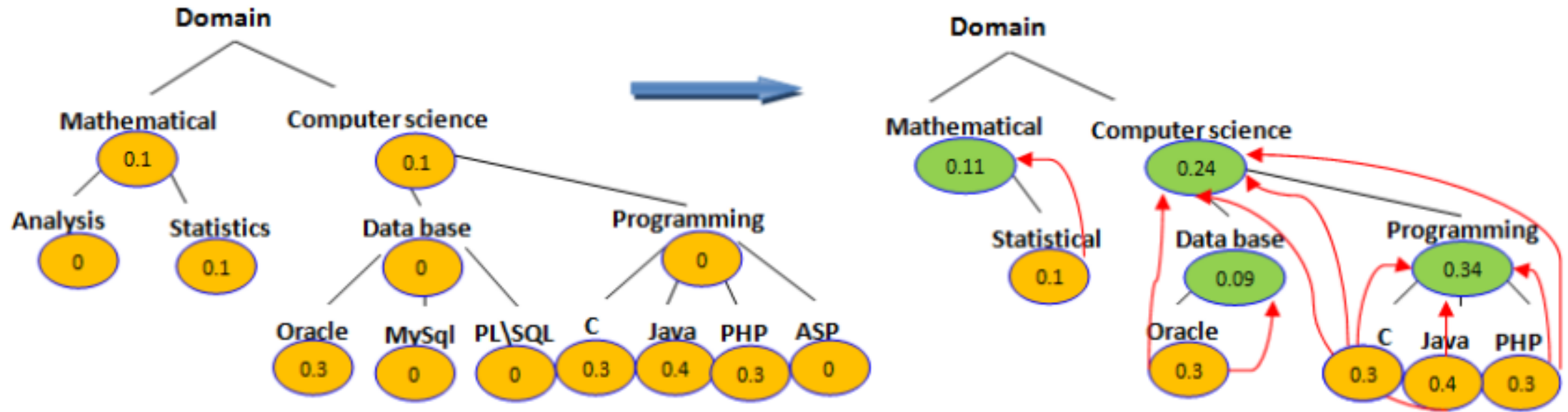
▶ Step 2

- ▶ For each leaf node which has a weight not null, the weights of his ancestors are computed by the formula below:

$$Weight(n_k, n_{fi}) = Weight(n_k) + Weight(n_{fi})^{distance(n_k, n_{fi}) + 1} \quad (2)$$



Relevance Propagation



more important

Result: $\overrightarrow{d_{Dim_z}} = \text{computer-science}(0.24), \underline{\text{programming}}(0.34), \text{Java}(0.4), \text{php}(0.3), \text{c}(0.3), \underline{\text{database}}(0.09), \text{oracle}(0.3), \text{math}(0.11), \text{statistics}(0.1).$

Text-OLAP based IR

$$Sim(d, q) = \cos a = \frac{\sum_i w_{ti} * w_{qi}}{\sqrt{\sum_i w_{ti}^2 * \sum_i w_{qi}^2}}$$

$$Sim(d, q) = \sum_{i=1}^n (\alpha_i \times Sim(d_{Dim_i}, q_i))$$

$$Sim(d, p) = \sum_{i=1}^n (\alpha_i \times Sim(d_{Dim_i}, p_i))$$

$$ORank(d) = \beta \times Sim(d, q) + \gamma \times Sim(d, p)$$



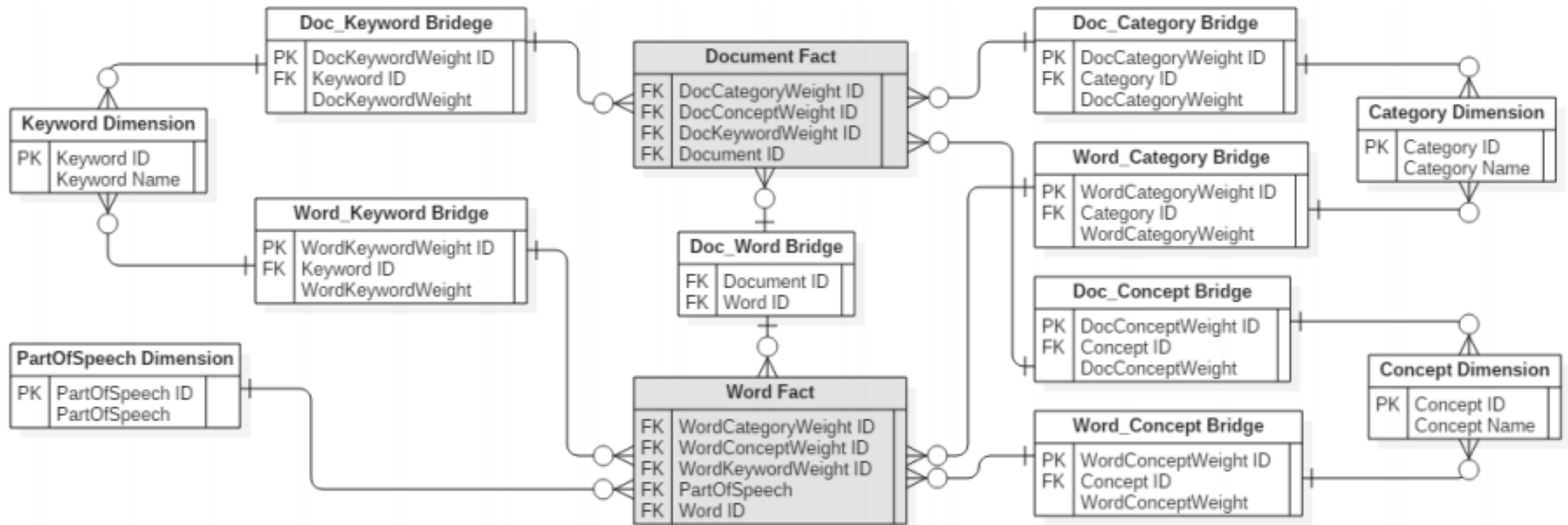
CXT-Cube



- ▶ For CV data
 - ▶ Fact table: CV documents
 - ▶ Dimension tables
 - ▶ Semantic dimensions: TOPIC, LOCATION
 - ▶ TOPIC : hierarchy of domain skills, ex) data science – computer science - science
 - ▶ LOCATION : ex) city-region-country



A Star Schema for Text Warehouse



DW-based Machine Learning

