# Multidimensional Mining of Massive Text Data

# Synthesis Lectures on Data Mining and Knowledge Discovery

Multidimensional Mining of Massive Text Data

Chao Zhang and Jiawei Han

# Multidimensional Mining of Massive Text Data

Chao Zhang
Georgia Institute of Technology

Jiawei Han
University of Illinois at Urbana-Champaign

## ABSTRACT

Unstructured text, as one of the most important data forms, plays a crucial role in data-driven decision making in domains ranging from social networking and information retrieval to scientific research and healthcare informatics. In many emerging applications, people's information need from text data is becoming multidimensional—they demand useful insights along multiple aspects from a text corpus. However, acquiring such multidimensional knowledge from massive text data remains a challenging task.

This book presents data mining techniques that turn unstructured text data into multidimensional knowledge. We investigate two core questions. (1) How does one identify task-relevant text data with declarative queries in multiple dimensions? (2) How does one distill knowledge from text data in a multidimensional space? To address the above questions, we develop a text cube framework. First, we develop a cube construction module that organizes unstructured data into a cube structure, by discovering latent multidimensional and multi-granular structure from the unstructured text corpus and allocating documents into the structure. Second, we develop a cube exploitation module that models multiple dimensions in the cube space, thereby distilling from user-selected data multidimensional knowledge. Together, these two modules constitute an integrated pipeline: leveraging the cube structure, users can perform multidimensional, multigranular data selection with declarative queries; and with cube exploitation algorithms, users can extract multidimensional patterns from the selected data for decision making.

The proposed framework has two distinctive advantages when turning text data into multidimensional knowledge: flexibility and label-efficiency. First, it enables acquiring multidimensional knowledge flexibly, as the cube structure allows users to easily identify task-relevant data along multiple dimensions at varied granularities and further distill multidimensional knowledge. Second, the algorithms for cube construction and exploitation require little supervision; this makes the framework appealing for many applications where labeled data are expensive to obtain.

## KEYWORDS

text mining, multidimensional analysis, data cube, limited supervision

# Contents

**3 Term-Level Taxonomy Generation** . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . **31**

*Jiaming Shen*
*University of Illinois at Urbana–Champaign*

**4 Weakly Supervised Text Classification** . . . . . . . . . . . . . . . . . . . . . . . . . . . **49**

*Yu Meng*
*University of Illinois at Urbana–Champaign*

**5**    **Weakly Supervised Hierarchical Text Classification**

*Yu Meng*

*University of Illinois at Urbana–Champaign*