# Mining
# of
# Massive
# Datasets

Anand Rajaraman

Kosmix, Inc.

Jeffrey D. Ullman

Stanford Univ.

# Preface

This book evolved from material developed over several years by Anand Raja-raman and Jeff Ullman for a one-quarter course at Stanford. The course CS345A, titled "Web Mining," was designed as an advanced graduate course, although it has become accessible and interesting to advanced undergraduates.

## What the Book Is About

At the highest level of description, this book is about data mining. However, it focuses on data mining of very large amounts of data, that is, data so large it does not fit in main memory. Because of the emphasis on size, many of our examples are about the Web or data derived from the Web. Further, the book takes an algorithmic point of view: data mining is about applying algorithms to data, rather than using data to "train" a machine-learning engine of some sort. The principal topics covered are:

1. Distributed file systems and map-reduce as a tool for creating parallel algorithms that succeed on very large amounts of data.

2. Similarity search, including the key techniques of minhashing and locality-sensitive hashing.

3. Data-stream processing and specialized algorithms for dealing with data that arrives so fast it must be processed immediately or lost.

4. The technology of search engines, including Google's PageRank, link-spam detection, and the hubs-and-authorities approach.

5. Frequent-itemset mining, including association rules, market-baskets, the A-Priori Algorithm and its improvements.

6. Algorithms for clustering very large, high-dimensional datasets.

7. Two key problems for Web applications: managing advertising and rec-ommendation systems.

## Prerequisites

CS345A, although its number indicates an advanced graduate course, has been found accessible by advanced undergraduates and beginning masters students. In the future, it is likely that the course will be given a mezzanine-level number. The prerequisites for CS345A are:

1. The first course in database systems, covering application programming in SQL and other database-related languages such as XQuery.

2. A sophomore-level course in data structures, algorithms, and discrete math.

3. A sophomore-level course in software systems, software engineering, and programming languages.

## Exercises

The book contains extensive exercises, with some for almost every section. We indicate harder exercises or parts of exercises with an exclamation point. The hardest exercises have a double exclamation point.

## Support on the Web

You can find materials from past offerings of CS345A at:

     http://infolab.stanford.edu/~ullman/mining/mining.html

There, you will find slides, homework assignments, project requirements, and in some cases, exams.

## Acknowledgements

Cover art is by Scott Ullman. We would like to thank Foto Afrati and Arun Marathe for critical readings of the draft of this manuscript. Errors were also reported by Apoorv Agarwal, Susan Biancani, Leland Chen, Shrey Gupta, Xie Ke, Haewoon Kwak, Ellis Lau, Ethan Lozano, Justin Meyer, Brad Penoff, Philips Kokoh Prasetyo, Angad Singh, Sandeep Sripada, Dennis Sidharta, Mark Storus, Roshan Sumbaly, and Tim Triche Jr. The remaining errors are ours, of course.

> A. R.
> J. D. U.
> Palo Alto, CA
> June, 2011

# Contents