



Statistical Methods for Recommender Systems

DEEPAK K. AGARWAL
BEE-CHUNG CHEN

Statistical Methods for Recommender Systems

DEEPAK K. AGARWAL

LinkedIn Corporation

BEE-CHUNG CHEN

LinkedIn Corporation



CAMBRIDGE
UNIVERSITY PRESS

CAMBRIDGE
UNIVERSITY PRESS

32 Avenue of the Americas, New York, NY 10013-2473, USA

Cambridge University Press is part of the University of Cambridge.

It furthers the University's mission by disseminating knowledge in the pursuit of education, learning, and research at the highest international levels of excellence.

www.cambridge.org

Information on this title: www.cambridge.org/9781107036079

© Deepak K. Agarwal and Bee-Chung Chen 2016

This publication is in copyright. Subject to statutory exception and to the provisions of relevant collective licensing agreements, no reproduction of any part may take place without the written permission of Cambridge University Press.

First published 2016

Printed in the United Kingdom by Clays, St Ives plc.

A catalog record for this publication is available from the British Library.

Library of Congress Cataloging in Publication Data

Agarwal, Deepak K., 1973– author.

Statistical methods for recommender systems / Deepak K. Agarwal, Yahoo!

Research, Bee Chung-Chen, Yahoo! Research.

pages cm

ISBN 978-1-107-03607-9

1. Recommender systems (Information filtering) – Statistical methods. 2. Expert systems (Computer science) – Statistical methods. I. Chung-Chen, Bee, author. II. Title.

QA76.76.E95A395 2016

006.3'3–dc23 2015026092

ISBN 978-1-107-03607-9 Hardback

Additional resources for this publication at

<https://github.com/bee chung/Latent-Factor-Models>

Cambridge University Press has no responsibility for the persistence or accuracy of URLs for external or third-party Internet websites referred to in this publication and does not guarantee that any content on such websites is, or will remain, accurate or appropriate.

Statistical Methods for Recommender Systems

Designing algorithms to recommend items such as news articles and movies to users is a challenging task in numerous web applications. The crux of the problem is to rank items based on past user responses to optimize for multiple objectives. Major technical challenges are high-dimensional prediction with sparse data and constructing high-dimensional sequential designs to collect data for user modeling and system design.

This comprehensive treatment of the statistical issues that arise in recommender systems includes detailed, in-depth discussions of current state-of-the-art methods such as adaptive sequential designs (multiarmed bandit methods), bilinear random-effects models (matrix factorization), and scalable model fitting using modern computing paradigms such as MapReduce. The authors draw on their vast experience working with such large-scale systems at Yahoo! and LinkedIn and bridge the gap between theory and practice by illustrating complex concepts with examples from applications with which they are directly involved.

DR. DEEPAK K. AGARWAL is a big data analyst with several years of experience developing and deploying state-of-the-art machine learning and statistical methods for improving the relevance of web applications. He is also experienced in conducting new scientific research to solve difficult big data problems, especially in the areas of recommender systems and computational advertising. He is a Fellow of the American Statistical Association and associate editor of top-tier journals in statistics.

DR. BEE-CHUNG CHEN is a leading technologist with extensive industrial and research experience in developing state-of-the-art recommender systems. He has been a key designer of the recommendation algorithms that power the LinkedIn home page and mobile feeds, the Yahoo! home page, Yahoo! News, and other sites. His research areas include recommender systems, data mining, machine learning, and big data analytics.

For Bharati Agarwal and Shiao-Ching Chung

Preface

What This Book Is About

Recommender systems are automated computer programs that match items to users in different contexts. Such systems are ubiquitous and have become an integral part of our daily lives. Examples include recommending products to users on a site like Amazon, recommending content to users visiting a website like Yahoo!, recommending movies to users on a site like Netflix, recommending jobs to users on a site like LinkedIn, and so on. The matching algorithms are constructed using large amounts of high-frequency data obtained from past user interactions with items. The algorithms are statistical in nature and involve challenges in areas like sequential decision processes, modeling interactions with very high-dimensional categorical data, and developing scalable statistical methods. New methodologies in this area require close collaboration among computer scientists, machine learners, statisticians, optimization experts, system experts, and, of course, domain experts. It is one of the most exciting applications of big data.

Why We Wrote This Book

Although much has been written about recommender systems in various fields, such as computer science, machine learning, and statistics, focusing on specific aspects of the problem, a comprehensive treatment of all statistical issues and how they are interrelated is lacking. We came to this realization while deploying such systems at Yahoo! and LinkedIn. For instance, much of the focus in statistics and machine learning is on building models that minimize out-of-sample predictive error. However, this does not address all aspects of practical importance. Statistically, a recommender system is a high-dimensional sequential process, and it is equally important to study issues like design of

experiments as it is to develop sophisticated statistical models. In fact, the two are closely related – efficient design needs models to tame the curse of dimensionality. Also, most existing work in the literature tends to build models for univariate response, such as movie ratings, purchases, and click rates. With the advent of social media outlets like Facebook, LinkedIn, and Twitter, multiple responses are available. For instance, one may want to model click rates, share rates, and tweet rates simultaneously for a news recommender application. Such multivariate response models are challenging to build. Finally, given the machinery to obtain such multivariate predictions, how does one construct utility functions to make recommendations? Is it more important to optimize share rates relative to click rates? Answers to these types of questions can be obtained through multiobjective optimization working in close collaboration with domain experts to elicit some utility parameters.

The goal of this book is to provide a comprehensive discussion of all such issues that arise in the context of recommender systems. This is in addition to a detailed and in-depth discussion of current state-of-the-art statistical methods that include techniques like adaptive sequential designs (multiarmed bandit methods), bilinear random-effects models (matrix factorization), and scalable model fitting using modern-distributed computing infrastructure. Our goal in writing this book is to draw on our vast experience working with such large-scale systems in industrial settings and to bring these issues to the attention of the statistical, machine learning, and computer science communities. We believe this will be beneficial in a number of ways. It may help in advancing methodological research in high-dimensional and big data statistics, especially for web applications. We understand that conducting such research in an academic setting requires access to software that can run on massive data. To facilitate this, we supplement the book with open source software: <https://github.com/bee chung/Latent-Factor-Models>. We also believe the book will help in bridging the gap between theory and applications. It will provide problem owners with a good understanding of the statistical issues involved and modelers with an in-depth understanding of statistical issues that arise in practical applications that are rather complex.

Organization

We divide the content of the book into three parts.

In Part I, we introduce the recommender system problem, challenges in the problem, main ideas used to tackle the challenges, and the required background knowledge. In Chapter 2, we give an overview of classical methods

that have been used to develop recommender systems. Such methods involve characterizing users and items as feature vectors and then scoring user-item pairs based on some similarity function, standard supervised learning, or collaborative filtering. These classical methods usually ignore the explore-exploit trade-off in recommender problems. Hence, in Chapter 3, we discuss the importance of this issue and introduce the main ideas that will be used to solve the issue in later chapters. Before we delve into technical solutions, in Chapter 4, we review a variety of methods for evaluating the performance of different recommendation algorithms.

In Part II, we provide detailed solutions to common problem settings. We start with an introduction to various problem settings and an example system architecture in Chapter 5, and then we devote the next three chapters to three common problem settings. Chapter 6 provides solutions to the most-popular recommendation problem, with a special focus on the explore-exploit aspect. Chapter 7 deals with personalized recommendation through feature-based regression, with an emphasis on how to continuously update the model(s) to leverage the most recent user-item interaction data and quickly converge to a good solution. Chapter 8 extends the methods developed in Chapter 7 from feature-based regression to factor models (matrix factorization) and, at the same time, provides a natural solution to the cold-start problem in factor models.

In Part III, we present three advanced topics. In Chapter 9, we present a factorization model that simultaneously identifies topics in items and users' affinities with different topics through a modified matrix factorization model that uses the latent Dirichlet allocation (LDA) topic model. In Chapter 10, we investigate context-dependent recommender problems, in which the recommended items not only need to have high affinity with the user but also have to be relevant to the context (e.g., recommending items related to a news article that the user is currently reading). In Chapter 11, we discuss a principled framework for optimizing multiple objectives based on a constrained optimization approach, where we seek to maximize one objective (e.g., revenue) subject to bounded loss in other objectives (e.g., no more than 5 percent loss in clicks).

Limitations

Like all books, ours has limitations. We do not provide an in-depth coverage of modern computational paradigms, such as Spark, that can be used to fit some of the models presented at scale. Online evaluation of models when users form a social graph cannot be done properly with traditional experimental design methods. New techniques that can adjust for interference because of social

graphs need to be developed. We do not cover such advanced topics in this book. Throughout, we address the problem of recommendations through a response prediction approach using regression as our main tool. This is primarily because we believe that output from these models is easy to combine with downstream utilities. We do not provide a comprehensive coverage of methods that are based on direct optimization of ranking loss functions. A comparison of the two approaches would also be a worthwhile topic for discussion.

Acknowledgement

Our special thanks to Raghu Ramakrishnan, Liang Zhang, Xuanhui Wang, Pradheep Elango, Bo Long, Bo Pang, Rajiv Khanna, Nitin Motgi, Seung-Taek Park, Scott Roy, Joe Zachariah for many insightful discussions and collaboration. We would also like to thank our colleagues both at Yahoo! and LinkedIn for all the encouragement and support without which many of the ideas we had would not see the light of the day.

Contents

Preface *page ix*

PART I INTRODUCTION

1	Introduction	3
	1.1 Overview of Recommender Systems for Web Applications	4
	1.2 A Simple Scoring Model: Most-Popular Recommendation	10
	Exercises	14
2	Classical Methods	15
	2.1 Item Characterization	16
	2.2 User Characterization	23
	2.3 Feature-Based Methods	25
	2.4 Collaborative Filtering	31
	2.5 Hybrid Methods	36
	2.6 Summary	37
	Exercises	38
3	Explore-Exploit for Recommender Problems	39
	3.1 Introduction to the Explore-Exploit Trade-off	40
	3.2 Multiarmed Bandit Problem	41
	3.3 Explore-Exploit in Recommender Systems	48
	3.4 Explore-Exploit with Data Sparsity	50
	3.5 Summary	54
	Exercise	54

4	Evaluation Methods	55
4.1	Traditional Offline Evaluation	56
4.2	Online Bucket Tests	66
4.3	Offline Simulation	70
4.4	Offline Replay	73
4.5	Summary	77
	Exercise	78
 PART II COMMON PROBLEM SETTINGS		
5	Problem Settings and System Architecture	81
5.1	Problem Settings	81
5.2	System Architecture	89
6	Most-Popular Recommendation	94
6.1	Example Application: Yahoo! Today Module	95
6.2	Problem Definition	96
6.3	Bayesian Solution	98
6.4	Non-Bayesian Solutions	107
6.5	Empirical Evaluation	109
6.6	Large Content Pools	117
6.7	Summary	118
	Exercises	119
7	Personalization through Feature-Based Regression	120
7.1	Fast Online Bilinear Factor Model	122
7.2	Offline Training	126
7.3	Online Learning	131
7.4	Illustration on Yahoo! Data Sets	134
7.5	Summary	141
	Exercise	141
8	Personalization through Factor Models	142
8.1	Regression-Based Latent Factor Model (RLFM)	142
8.2	Fitting Algorithms	150
8.3	Illustration of Cold Start	164
8.4	Large-Scale Recommendation of Time-Sensitive Items	167
8.5	Illustration of Large-Scale Problems	172
8.6	Summary	182
	Exercise	182

PART III ADVANCED TOPICS

9	Factorization through Latent Dirichlet Allocation	185
9.1	Introduction	185
9.2	Model	186
9.3	Training and Prediction	191
9.4	Experiments	198
9.5	Related Work	203
9.6	Summary	204
10	Context-Dependent Recommendation	206
10.1	Tensor Factorization Models	207
10.2	Hierarchical Shrinkage	211
10.3	Multifaceted News Article Recommendation	218
10.4	Related-Item Recommendation	233
10.5	Summary	235
11	Multiobjective Optimization	237
11.1	Application Setting	238
11.2	Segmented Approach	239
11.3	Personalized Approach	243
11.4	Approximation Methods	248
11.5	Experiments	250
11.6	Related Work	261
11.7	Summary	262
	<i>Endnotes</i>	263
	<i>References</i>	265
	<i>Index</i>	273

PART I

Introduction

1

Introduction

Recommender systems (or recommendation systems) are computer programs that recommend the “best” items to users in different contexts. The notion of a best match is typically obtained by optimizing for objectives like total clicks, total revenue, and total sales. Such systems are ubiquitous on the web and form an integral part of our daily lives. Examples include product recommendations to users on an e-commerce site to maximize sales; content recommendations to users visiting a news site to maximize total clicks; movie recommendations to maximize user engagement and increase subscriptions; or job recommendations on a professional network site to maximize job applications. Input to these algorithms typically consists of information about users, items, contexts, and feedback that is obtained when users interact with items.

Figure 1.1 shows an example of a typical web application that is powered by a recommender system. A user uses a web browser to visit a web page. The browser then submits an HTTP request to the web server that hosts the page. To serve recommendations on the page (e.g., popular news stories on a news portal page), the web server makes a call to a recommendation service that retrieves a set of items and renders them on the web page. Such a service typically performs a large number of different types of computations to select the best items. These computations are often a hybrid of both offline and real-time computations, but they must adhere to strict efficiency requirements to ensure quick page load time (typically hundreds of milliseconds). Once the page loads, the user may interact with items through actions like clicks, likes, or shares. Data obtained through such interactions provide a feedback loop to update the parameters of the underlying recommendation algorithm and to improve the performance of the algorithm for future user visits. The frequency of such parameter updates depends on the application. For instance, if items are time sensitive or ephemeral, as in the case of news recommendations, parameter updates must be done frequently (e.g., every few minutes). For

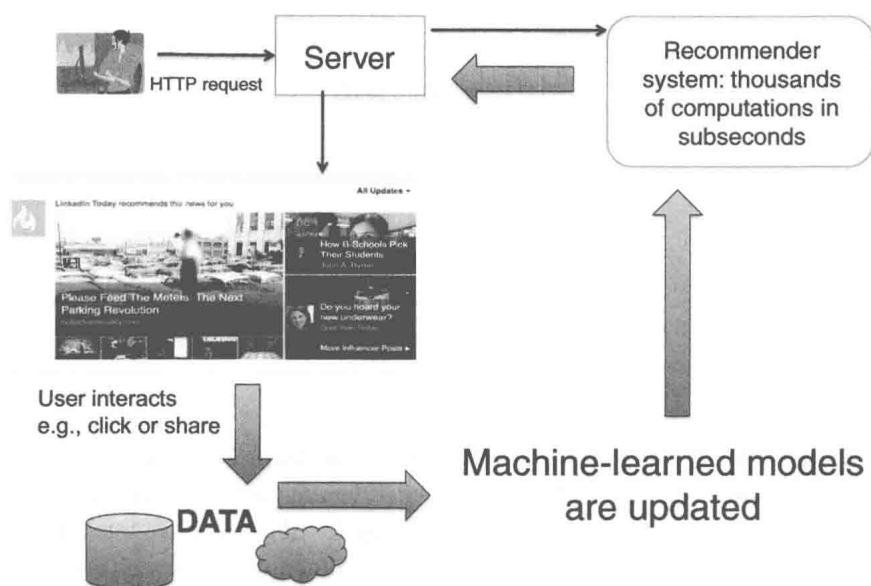


Figure 1.1. A typical recommender system.

other applications where items have a relatively longer lifetime (e.g., movie recommendations), parameter updates can happen less frequently (e.g., daily) without significant degradation in overall performance.

Algorithms that facilitate selection of best items are crucial to the success of recommender systems. This book provides a comprehensive description of statistical and machine learning methods on which we believe such algorithms should be based. For the sake of simplicity, we loosely refer to these algorithms as recommender systems throughout this book, but note that they only represent one component (albeit a crucial one) of the end-to-end process required to serve items to users in a scalable fashion.

1.1 Overview of Recommender Systems for Web Applications

Before developing a recommender system, it is important to consider the following questions.

- *What input signals are available?* When building machine-learned models of what items a user is likely to interact with in a given context, we can draw on many signals, including the content and source of each item; a