# Statistical Foundations of Data Science

Jianqing Fan

Runze Li

Cun-Hui Zhang

Hui Zou

# Preface

Big data are ubiquitous. They come in varying volume, velocity, and variety. They have a deep impact on systems such as storages, communications and computing architectures and analysis such as statistics, computation, optimization, and privacy. Engulfed by a multitude of applications, data science aims to address the large-scale challenges of data analysis, turning big data into smart data for decision making and knowledge discoveries. Data science integrates theories and methods from statistics, optimization, mathematical science, computer science, and information science to extract knowledge, make decisions, discover new insights, and reveal new phenomena from data. The concept of data science has appeared in the literature for several decades and has been interpreted differently by different researchers. It has nowadays become a multi-disciplinary field that distills knowledge in various disciplines to develop new methods, processes, algorithms and systems for knowledge discovery from various kinds of data, which can be either low or high dimensional, and either structured, unstructured or semi-structured. Statistical modeling plays critical roles in the analysis of complex and heterogeneous data and quantifies uncertainties of scientific hypotheses and statistical results.

This book introduces commonly-used statistical models, contemporary statistical machine learning techniques and algorithms, along with their mathematical insights and statistical theories. It aims to serve as a graduate-level textbook on the statistical foundations of data science as well as a research monograph on sparsity, covariance learning, machine learning and statistical inference. For a one-semester graduate level course, it may cover Chapters 2, 3, 9, 10, 12, 13 and some topics selected from the remaining chapters. This gives a comprehensive view on statistical machine learning models, theories and methods. Alternatively, one-semester graduate course may cover Chapters 2, 3, 5, 7, 8 and selected topics from the remaining chapters. This track focuses more on high-dimensional statistics, model selection and inferences but both paths emphasize a great deal on sparsity and variable selections.

Frontiers of scientific research rely on the collection and processing of massive complex data. Information and technology allow us to collect big data of unprecedented size and complexity. Accompanying big data is the rise of dimensionality and high dimensionality characterizes many contemporary statistical problems, from sciences and engineering to social science and humanities. Many traditional statistical procedures for finite or low-dimensional data are still useful in data science, but they become infeasible or ineffective for dealing with high-dimensional data. Hence, new statistical methods are indispensable. The authors have worked on high-dimensional statistics for two decades, and started to write the book on the topics of high-dimensional data analysis over a decade ago. Over the last decide, there have been surges in interest and exciting developments in high-dimensional and big data. This led us to concentrate mainly on statistical aspects of data science.

We aim to introduce commonly-used statistical models, methods and pro-

cedures in data science and provide readers with sufficient and sound theoretical justifications. It has been a challenge for us to balance statistical theories and methods and to choose the topics and works to cover since the amount of publications in this emerging area is enormous. Thus, we focus on the foundational aspects that are related to sparsity, covariance learning, machine learning, and statistical inference.

Sparsity is a common assumption in the analysis of high-dimensional data. By sparsity, we mean that only a handful of features embedded in a huge pool suffice for certain scientific questions or predictions. This book introduces various regularization methods to deal with sparsity, including how to determine penalties and how to choose tuning parameters in regularization methods and numerical optimization algorithms for various statistical models. They can be found in Chapters 3–6 and 8.

High-dimensional measurements are frequently dependent, since these variables often measure similar things, such as aspects of economics or personal health. Many of these variables have heavy tails due to big number of collected variables. To model the dependence, factor models are frequently employed, which exhibit low-rank plus sparse structures in data matrices and can be solved by robust principal component analysis from high-dimensional covariance. Robust covariance learning, principal component analysis, as well as their applications to community detection, topic modeling, recommender systems, ect. are also a feature of this book. They can be found in Chapters 9–11. Note that factor learning or more generally latent structure learning can also be regarded as unsupervised statistical machine learning.

Machine learning is critical in analyzing high-dimensional and complex data. This book also provides readers with a comprehensive account on statistical machine learning methods and algorithms in data science. We introduce statistical procedures for supervised learning in which the response variable (often categorical) is available and the goal is to predict the response based on input variables. This book also provides readers with statistical procedures for unsupervised learning, in which the responsible variable is missing and the goal concentrates on learning the association and patterns among a set of input variables. Feature creations and sparsity learning also arise in these problems. See Chapters 2, 12–14 for details.

Statistical inferences on high-dimensional data are another focus of this book. Statistical inferences require one to characterize the uncertainty, estimate the standard errors of the estimated parameters of primary interest and derive the asymptotic distributions of the resulting estimates. This is very challenging under the high-dimensional regime. See Chapter 7.

Fueled by the surging demands on processing high-dimensional and big data, there have been rapid and vast developments in high-dimensional statistics and machine learning over the last decade, contributed by data scientists from various fields such as statistics, computer science, information theory, applied and computational mathematics, among others. Even though we have narrowed the scope of the book to the statistical aspects of data science, the

field is still too broad for us to cover. Many important contributions that do not fit our presentation have been omitted. Conscientious effort was made in the composition of the reference list and bibliographical notes, but they merely reflect our immediate interests. Omissions and discrepancies are inevitable. We apologize for their occurrence.

Although we all contribute to various chapters and share the responsibility for the whole book, Jianqing Fan was the lead author for Chapters 1, 3 and 9–11, 14 and some sections in other chapters, Runze Li for Chapters 5, and 8 and part of Chapters 6–7, Cun-Hui Zhang for Chapters 4 and 7, and Hui Zou for Chapters 2, 6, 11 and 12 and part of Chapter 5. In addition, Jianqing Fan and Runze Li oversaw the whole book project.

Many people have contributed importantly to the completion of this book. In particular, we would like to thank the editor, John Kimmel, who has been extremely helpful and patient with us for over 10 years! We greatly appreciate a set of around 10 anonymous reviewers for valuable comments that lead to the improvement of the book. We are particularly grateful to Cong Ma and Yiqiao Zhong for preparing a draft of Chapter 14, to Zhao Chen for helping us with putting our unsorted and non-uniform references into the present form, to Tracy Ke, Bryan Kelly, Dacheng Xiu and Jia Wang for helping us with constructing Figure 1.3, and to Boxiang Wang, Yi Yang for helping produce some figures in Chapter 12. Various people have carefully proof-read certain chapters of the book and made useful suggestions. They include Krishna Balasubramanian, Pierre Bayle, Elynn Chen, Wenyan Gong, Yongyi Guo, Cong Ma, Igor Silin, Qiang Sun, Francesca Tang, Bingyan Wang, Kaizheng Wang, Weichen Wang, Yuling Yan, Zhuoran Yang, Mengxin Yu, Wenxin Zhou, Yifeng Zhou, and Ziwei Zhu. We owe them many thanks.

In the spring semester of 2019, we used a draft of this book as a textbook for a first-year graduate course at Princeton University and a senior graduate topic course at the Pennsylvania State University. We would like to thank the graduate students in the classes for their careful readings. In particular, we are indebted to Cong Ma, Kaizheng Wang and Zongjun Tan for assisting in preparing the homework problems at Princeton, most of which are now a part of our exercise at the end of each chapter. At Princeton, we covered chapters 2-3, 5, 8.1, 8.3, 9–14.

We are very grateful for grant supports from National Science Foundation and National Institutes of Health on our research. Finally, we would like to thank our families and our parents for their love and support.

<div style="text-align: right">

Jianqing Fan
Runze Li
Cun-Hui Zhang
Hui Zou

January 2020.

</div>

# Contents