

ROAD TO STATISTICAL BIOINFORMATICS

Jae K. Lee

*Department of Public Health Science, University of Virginia,
Charlottesville, Virginia, USA*

There has been a great explosion of biological data and information in recent years, largely due to the advances of various high-throughput biotechnologies such as mass spectrometry, high throughput sequencing, and many genome-wide SNP profiling, RNA gene expression microarray, protein mass spectrometry, and many other recent high-throughput biotechniques (Weinstein et al., 2002). Furthermore, powerful computing systems and fast Internet connections to large worldwide biological databases enable individual laboratory researchers to easily access an unprecedentedly huge amount of biological data. Such enormous data are often too overwhelming to understand and extract the most relevant information to each researcher's investigation goals. In fact, these large biological data are information rich and often contain much more information than the researchers who have generated such data may have anticipated. This is why many major biomedical research institutes have made significant efforts to freely share such data with general public researchers. Bioinformatics is the emerging science field concerned with the development of various analysis methods and tools for investigating such large biological data efficiently and rigorously. This kind of development requires many different components: powerful computer systems to archive and process such data, effective database designs to extract and integrate information from various heterogeneous biological databases, and efficient analysis techniques to investigate and analyze these large databases. In particular, analysis of these massive biological data is extremely challenging for the following reasons.

CHALLENGE 1: MULTIPLE-COMPARISONS ISSUE

Analysis techniques on high-throughput biological data are required to carefully handle and investigate an astronomical number of candidate targets and possible mechanisms, most of which are false positives, from such massive data (Tusher

et al., 2001). For example, a traditional statistical testing criterion which allows 5% false-positive error (or significance level) would identify ~ 500 false positives from 10K microarray data between two biological conditions of interest even though no real biologically differentially regulated genes exist between the two. If a small number of, for example, 100, genes that are actually differentially regulated exist, such real differential expression patterns will be mixed with the above 500 false positives without any a priori information to discriminate the true positives from the false positives. Then, confidence on the 600 targets that were identified by such a statistical testing may not be high. Simply tightening such a statistical criterion will result in a high false-negative error rate, without being able to identify many important real biological targets. This kind of pitfall, the so-called *multiple-comparisons issue*, becomes even more serious when biological mechanisms such as certain signal transduction and regulation pathways that involve multiple targets are searched from such biological data; the number of candidate pathway mechanisms to be searched grows exponentially, for example, $10!$ for 10-gene sequential pathway mechanisms. Thus, no matter how powerful a computer system can handle a given computational task, it is prohibitive to tackle such problems by exhaustive computational search and comparison for these kinds of problems. Many current biological problems have been theoretically proven to be NP (nonpolynomial) hard in computer science, implying that no finite (polynomial) computational algorithm can search all possible solutions as the number of biological targets involved in such a solution becomes too large. More importantly, this kind of exhaustive search is simply prone to the risk of discovering numerous false positives. In fact, this is one of the most difficult challenges in investigating current large biological databases and is why only heuristic algorithms that tightly control such a high false positive error rate and investigate a very small portion of all possible solutions are often sought for many biological problems. Thus, the success of many bioinformatics studies critically depends on the construction and use of effective and efficient heuristic algorithms, most of which are based on probabilistic modeling and statistical inference techniques that can maximize the statistical power of identifying true positives while rigorously controlling their false positive error rates.

CHALLENGE 2: HIGH-DIMENSIONAL BIOLOGICAL DATA

The second challenge is the high-dimensional nature of biological data in many bioinformatics studies. When biological data are simultaneously generated with many gene targets, their data points become dramatically sparse in the corresponding high-dimensional data space. It is well known that mathematical and computational approaches often fail to capture such high-dimensional phenomena accurately (Tamayo et al., 1999). For example, many statistical algorithms cannot easily move between local maxima in a high-dimensional space. Also, inference by combining several disjoint lower dimensional phenomena may not provide the correct understanding on the real phenomena in their joint, high-dimensional space. It is therefore important to understand statistical dimension reduction techniques that

can reduce high-dimensional data problems into lower dimensional ones while the important variation of interest in biological data is preserved.

CHALLENGE 3: SMALL- n AND LARGE- p PROBLEM

The third challenge is the so-called “small- n and large- p ” problem. Desired performance of conventional statistical methods is achieved when the sample size, namely n , of the data, the number of independent observations of event, is much larger than the number of parameters, say p , which need to be inferred by statistical inference (Jain et al., 2003). In many bioinformatics problems, this situation is often completely reversed. For example, in a microarray study, tens of thousands of gene transcripts’ expression patterns may become candidate prediction factors for a biological phenomenon of interest (e.g., tumor sensitivity vs. resistance to a chemotherapeutic compound) but the number of independent observations (e.g., different patient biopsy samples) is often at most a few tens or smaller. Due to the experimental costs and limited biological materials, the number of independent replicated samples can be sometimes extremely small, for example, two or three, or unavailable. In these cases, most traditional statistical approaches often perform very poorly. Thus, it is also important to select statistical analysis tools that can provide both high specificity and high sensitivity under these circumstances.

CHALLENGE 4: NOISY HIGH-THROUGHPUT BIOLOGICAL DATA

The fourth challenge is due to the fact that high-throughput biotechnical data and large biological databases are inevitably noisy because biological information and signals of interest are often observed with many other random or biased factors that may obscure main signals and information of interest (Cho and Lee, 2004). Therefore, investigations on large biological data cannot be successfully performed unless rigorous statistical algorithms are developed and effectively utilized to reduce and decompose various sources of error. Also, careful assessment and quality control of initial data sets is critical for all subsequent bioinformatics analyses.

CHALLENGE 5: INTEGRATION OF MULTIPLE, HETEROGENEOUS BIOLOGICAL DATA INFORMATION

The last challenge is the integration of information often from multiple heterogeneous biological and clinical data sets, such as large gene functional and annotation databases, biological subjects’ phenotypes, and patient clinical information. One of the main goals in performing high-throughput biological experiments is to identify the most important critical biological targets and mechanisms highly associated with biological subjects’ phenotypes, such as patients’ prognosis and therapeutic response

(Pittman et al., 2004). In these cases, multiple large heterogeneous datasets need to be combined in order to discover the most relevant molecular targets. This requires combining multiple datasets with very different data characteristics and formats, some of which cannot easily be integrated by standard statistical inference techniques, for example, the information from genomic and proteomic expression data and reported pathway mechanisms in the literature. It will be extremely important to develop and use efficient yet rigorous analysis tools for integrative inference on such complex biological data information beyond the individual researcher's manual and subjective integration.

In this book, we introduce the statistical concepts and techniques that can overcome these challenges in studying various large biological datasets. Researchers with biological or biomedical backgrounds may not be able, or may not need, to learn advanced mathematical and statistical techniques beyond the intuitive understanding of such topics for their practical applications. Thus, we have organized this book for life science researchers to efficiently learn the most relevant statistical concepts and techniques for their specific biological problems. We believe that this composition of the book will help nonstatistical researchers to minimize unnecessary efforts in learning statistical topics that are less relevant to their specific biological questions, yet help them learn and utilize rigorous statistical methods directly relevant to those problems. Thus, while this book can serve as a general reference for various concepts and methods in statistical bioinformatics, it is also designed to be effectively used as a textbook for a semester or shorter length course as below. In particular, the chapters are divided into four blocks of different statistical issues in analyzing large biological datasets (Fig. 1.1):

- I. *Statistical Foundation* Probability theories (Chapter 2), statistical quality control (Chapter 3), statistical tests (Chapter 4)

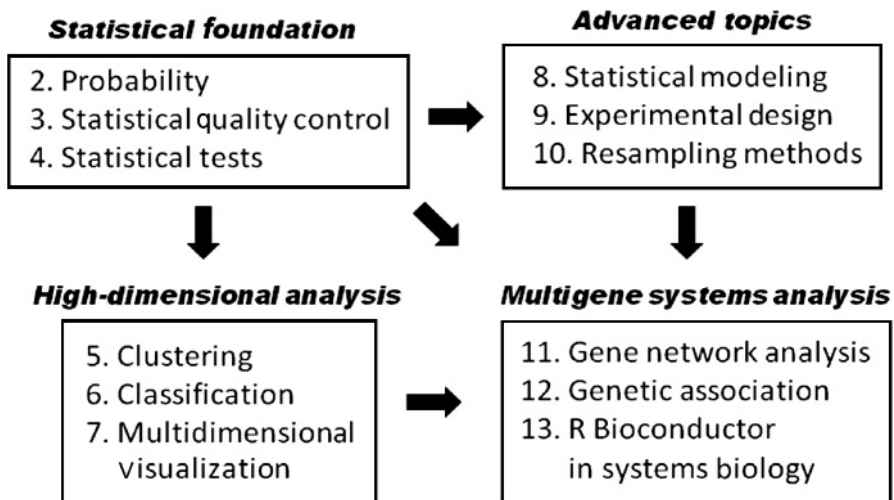


Figure 1.1 Possible course structure.

- II. *High-Dimensional Analysis* Clustering analysis (Chapter 5), classification analysis (Chapter 6), multidimensional visualization (Chapter 7)
- III. *Advanced Analysis Topics* Statistical modeling (Chapter 8), experimental design (Chapter 9), statistical resampling methods (Chapter 10)
- IV. *Multigene Analysis in Systems Biology* Genetic network analysis (Chapter 11), genetic association analysis (Chapter 12), R Bioconductor tools in systems biology (Chapter 13)

The first block of chapters will be important, especially for students who do not have a strong statistical background. These chapters will provide general backgrounds and terminologies to initiate rigorous statistical analysis on large biological datasets and to understand more advanced analysis topics later. Students with a good statistical understanding may also quickly review these chapters since there are certain key concepts and techniques (especially in Chapters 3 and 4) that are relatively new and specialized for analyzing large biological datasets.

The second block consists of analysis topics frequently used in investigating high-dimensional biological data. In particular, clustering and classification techniques, by far, are most commonly used in many practical applications of high-throughput data analysis. Various multidimensional visualization tools discussed in Chapter 7 will also be quite handy in such investigations.

The third block deals with more advanced topics in large biological data analysis, including advanced statistical modeling for complex biological problems, statistical resampling techniques that can be conveniently used with the combination of classification (Chapter 6) and statistical modeling (Chapter 8) techniques, and experimental design issues in high-throughput microarray studies.

The final block contains concise description of the analysis topics in several active research areas of multigene network and genetic association analysis as well as the R Bioconductor software in systems biology analysis. These will be quite useful for performing challenging gene network and multigene investigations in the fast-growing systems biology field.

These four blocks of chapters can be followed with the current order for a full semester-length course. However, except for the first block, the following three blocks are relatively independent of each other and can be covered (or skipped for specific needs and foci under a time constraint) in any order, as depicted in Figure 1.1. We hope that life science researchers who need to deal with challenging analysis issues in overwhelming large biological data in their specific investigations can effectively meet their learning goals in this way.

REFERENCES

- Cho, H., and Lee, J. K. (2004). Bayesian hierarchical error model for analysis of gene expression data. *Bioinformatics*, **20**(13): 2016–2025.
- Jain, N., et al. (2003). Local-pooled-error test for identifying differentially expressed genes with a small number of replicated microarrays. *Bioinformatics*, **19**(15): 1945–1951.

6 CHAPTER 1 ROAD TO STATISTICAL BIOINFORMATICS

- Pittman, J., et al. (2004). Integrated modeling of clinical and gene expression information for personalized prediction of disease outcomes. *Proc. Natl. Acad. Sci. U.S.A.*, **101**(22): 8431–8436.
- Tamayo, P., et al. (1999). Interpreting patterns of gene expression with self-organizing maps: Methods and application to hematopoietic differentiation. *Proc. Natl. Acad. Sci. U.S.A.*, **96**(6): 2907–2912.
- Tusher, V. G., Tibshirani, R., and Chu, G. (2001). Significance analysis of microarrays applied to the ionizing radiation response. *Proc. Natl. Acad. Sci. U.S.A.*, **98**(9): 5116–5121.
- Weinstein, J. N., et al. (2002). The bioinformatics of microarray gene expression profiling. *Cytometry*, **47**(1): 46–49.