

Chapter 1.15

Intelligent Data Analysis

Xiaohui Liu
Brunel University, UK

INTRODUCTION

Intelligent data analysis (IDA) is an interdisciplinary study concerned with the effective analysis of data. IDA draws the techniques from diverse fields, including artificial intelligence, databases, high-performance computing, pattern recognition, and statistics. These fields often complement each other (e.g., many statistical methods, particularly those for large data sets, rely on computation, but brute computing power is no substitute for statistical knowledge) (Berthold & Hand 2003; Liu, 1999).

BACKGROUND

The job of a data analyst typically involves problem formulation, advice on data collection (though it is not uncommon for the analyst to be asked to analyze data that have already been collected), effective data analysis, and interpretation and report of the finding. Data analysis is about the extraction of useful information from data and is

often performed by an iterative process in which exploratory analysis and confirmatory analysis are the two principal components.

Exploratory data analysis, or data exploration, resembles the job of a detective; that is, understanding evidence collected, looking for clues, applying relevant background knowledge, and pursuing and checking the possibilities that clues suggest.

Data exploration is not only useful for data understanding but also helpful in generating possibly interesting hypotheses for a later study—normally a more formal or confirmatory procedure for analyzing data. Such procedures often assume a potential model structure for the data and may involve estimating the model parameters and testing hypotheses about the model.

Over the last 15 years, we have witnessed two phenomena that have affected the work of modern data analysts more than any others. First, the size and variety of machine-readable data sets have increased dramatically, and the problem of data explosion has become apparent. Second, recent developments in computing have provided the

basic infrastructure for fast data access as well as many advanced computational methods for extracting information from large quantities of data. These developments have created a new range of problems and challenges for data analysts as well as new opportunities for intelligent systems in data analysis, and have led to the emergence of the field of intelligent data analysis (IDA), which draws the techniques from diverse fields, including artificial intelligence (AI), databases, high-performance computing, pattern recognition, and statistics. What distinguishes IDA is that it brings together often complementary methods from these diverse disciplines to solve challenging problems with which any individual discipline would find difficult to cope, and to explore the most appropriate strategies and practices for complex data analysis.

MAIN THRUST

In this paper, we will explore the main disciplines and associated techniques as well as applications to help clarify the meaning of intelligent data analysis, followed by a discussion of several key issues.

Statistics and Computing: Key Disciplines

IDA has its origins in many disciplines, principally statistics and computing. For many years, statisticians have studied the science of data analysis and have laid many of the important foundations. Many of the analysis methods and principles were established long before computers were born. Given that statistics are often regarded as a branch of mathematics, there has been an emphasis on mathematics rigor, a desire to establish that something is sensible on theoretical ground before trying it out on practical problems (Berthold & Hand, 2003). On the other hand, the computing community, particularly in machine

learning (Mitchell, 1997) and data mining (Wang, 2003) is much more willing to try something out (e.g., designing new algorithms) to see how they perform on real-world datasets, without worrying too much about the theory behind it.

Statistics is probably the oldest ancestor of IDA, but what kind of contributions has computing made to the subject? These may be classified into three categories. First, the basic computing infrastructure has been put in place during the last decade or so, which enables large-scale data analysis (e.g., advances in data warehousing and online analytic processing, computer networks, desktop technologies have made it possible to easily organize and move the data around for the analysis purpose). The modern computing processing power also has made it possible to efficiently implement some of the very computationally-intensive analysis methods such as statistical resampling, visualizations, large-scale simulation and neural networks, and stochastic search and optimization methods.

Second, there has been much work on extending traditional statistical and operational research methods to handle challenging problems arising from modern data sets. For example, in Bayesian networks (Ramoni et al., 2002), where the work is based on Bayesian statistics, one tries to make the ideas work on large-scale practical problems by making appropriate assumptions and developing computationally efficient algorithms; in support vector machines (Cristianini & Shawe-Taylor, 2000), where one tries to see how the statistical learning theory (Vapnik, 1998) could be utilized to handle very high-dimensional datasets in linear feature spaces; and in evolutionary computation (Eiben & Michalewicz, 1999) one tries to extend the traditional operational research search and optimization methods.

Third, new kinds of IDA algorithms have been proposed to respond to new challenges. Here are several examples of the novel methods with distinctive computing characteristics: powerful three-dimensional virtual reality visualization

5 more pages are available in the full version of this document, which may be purchased using the "Add to Cart" button on the publisher's webpage: www.igi-global.com/chapter/intelligent-data-analysis/24285

Related Content

A New Efficient and Effective Fuzzy Modeling Method for Binary Classification

T. Warren Liao (2013). *Contemporary Theory and Pragmatic Approaches in Fuzzy Computing Utilization* (pp. 41-59).

www.irma-international.org/chapter/new-efficient-effective-fuzzy-modeling/67481

An Optimal Configuration of Sensitive Parameters of PSO Applied to Textual Clustering

Reda Mohamed Hamou, Abdelmalek Amine, Mohamed Amine Boudiaand Ahmed Chaouki Lokbani (2019). *Exploring Critical Approaches of Evolutionary Computation* (pp. 196-214).

www.irma-international.org/chapter/an-optimal-configuration-of-sensitive-parameters-of-psy-applied-to-textual-clustering/208048

Performance Evaluation of Machine Learning for Recognizing Human Facial Emotions

Alti Adeland Ayeche Farid (2021). *International Journal of Intelligent Information Technologies* (pp. 1-17).

www.irma-international.org/article/performance-evaluation-of-machine-learning-for-recognizing-human-facial-emotions/286625

Big Data Analytics for Intrusion Detection: An Overview

Luis Filipe Diasand Miguel Correia (2020). *Handbook of Research on Machine and Deep Learning Applications for Cyber Security* (pp. 292-316).

www.irma-international.org/chapter/big-data-analytics-for-intrusion-detection/235047

On Multi-Fuzzy Rough Sets, Relations, and Topology

Gayathri Varmaand Sunil Jacob John (2019). *International Journal of Fuzzy System Applications* (pp. 101-119).

www.irma-international.org/article/on-multi-fuzzy-rough-sets-relations-and-topology/214942