

DATA SCIENCE WITH SEMANTIC TECHNOLOGIES

THEORY, PRACTICE AND APPLICATION

Edited by
Archana Patel
Narayan C. Debnath
Bharat Bhusan

 Scrivener
Publishing

WILEY

Table of Contents

[Cover](#)

[Title Page](#)

[Copyright](#)

[Preface](#)

[1 A Brief Introduction and Importance of Data Science](#)

[1.1 What is Data Science? What Does a Data Scientist Do?](#)

[1.2 Why Data Science is in Demand?](#)

[1.3 History of Data Science](#)

[1.4 How Does Data Science Differ from Business Intelligence?](#)

[1.5 Data Science Life Cycle](#)

[1.6 Data Science Components](#)

[1.7 Why Data Science is Important](#)

[1.8 Current Challenges](#)

[1.9 Tools Used for Data Science](#)

[1.10 Benefits and Applications of Data Science](#)

[1.11 Conclusion](#)

[References](#)

[2 Exploration of Tools for Data Science](#)

[2.1 Introduction](#)

[2.2 Top Ten Tools for Data Science](#)

[2.3 Python for Data Science](#)

[2.4 R Language for Data Science](#)

[2.5 SQL for Data Science](#)

[2.6 Microsoft Excel for Data Science](#)

[2.7 D3.JS for Data Science](#)

[2.8 Other Important Tools for Data Science](#)

[2.9 Conclusion](#)

[References](#)

[3 Data Modeling as Emerging Problems of Data Science](#)

[3.1 Introduction](#)

[3.2 Data](#)

[3.3 Data Model Design](#)

[3.4 Data Modeling](#)

[3.5 Polyglot Persistence Environment](#)

[References](#)

[4 Data Management as Emerging Problems of Data Science](#)

[4.1 Introduction](#)

[4.2 Perspective and Context](#)

[4.3 Data Distribution](#)

[4.4 CAP Theorem](#)

[4.5 Polyglot Persistence](#)

[References](#)

[5 Role of Data Science in Healthcare](#)

[5.1 Predictive Modeling—Disease Diagnosis and Prognosis](#)

[5.2 Preventive Medicine—Genetics/Molecular Sequencing](#)

[5.3 Personalized Medicine](#)

[5.4 Signature Biomarkers Discovery from High Throughput Data](#)

[Conclusion](#)

References

6 Partitioned Binary Search Trees (P(h)-BST): A Data Structure for Computer RAM

6.1 Introduction

6.2 P(h)-BST Structure

6.3 Maintenance Operations

6.4 Insert and Delete Algorithms

6.5 P(h)-BST as a Generator of Balanced Binary Search Trees

6.6 Simulation Results

6.7 Conclusion

Acknowledgments

References

7 Security Ontologies: An Investigation of Pitfall Rate

7.1 Introduction

7.2 Secure Data Management in the Semantic Web

7.3 Security Ontologies in a Nutshell

7.4 InFra_OE Framework

7.5 Conclusion

References

8 IoT-Based Fully-Automated Fire Control System

8.1 Introduction

8.2 Related Works

8.3 Proposed Architecture

8.4 Major Components

8.5 Hardware Interfacing

8.6 Software Implementation

8.7 Conclusion

References

9 Phrase Level-Based Sentiment Analysis Using Paired Inverted Index and Fuzzy Rule

9.1 Introduction

9.2 Literature Survey

9.3 Methodology

9.4 Conclusion

References

10 Semantic Technology Pillars: The Story So Far

10.1 The Road that Brought Us Here

10.2 What is a Semantic Pillar?

10.3 The Foundation Semantic Pillars: IRI's, RDF, and RDFS

10.4 The Semantic Upper Pillars: OWL, SWRL, SPARQL, and SHACL

10.5 Conclusion

References

11 Evaluating Richness of Security Ontologies for Semantic Web

11.1 Introduction

11.2 Ontology Evaluation: State-of-the-Art

11.3 Security Ontology

11.4 Richness of Security Ontologies

11.5 Conclusion

References

12 Health Data Science and Semantic Technologies

12.1 Health Data

12.2 Data Science

12.3 Health Data Science

12.4 Examples of Health Data Science Applications

[12.5 Health Data Science Challenges](#)

[12.6 Health Data Science and Semantic Technologies](#)

[12.7 Application of Data Science for COVID-19](#)

[12.8 Data Challenges During COVID-19 Outbreak](#)

[12.9 Biomedical Data Science](#)

[12.10 Conclusion](#)

[References](#)

[13 Hybrid Mixed Integer Optimization Method for Document Clustering Based on Semantic Data Matrix](#)

[13.1 Introduction](#)

[13.2 A Method for Constructing a Semantic Matrix of Relations Between Documents and Taxonomy Concepts](#)

[13.3 Mathematical Statements for Clustering Problem](#)

[13.4 Heuristic Hybrid Clustering Algorithm](#)

[13.5 Application of a Hybrid Optimization Algorithm for Document Clustering](#)

[13.6 Conclusion](#)

[References](#)

[14 Role of Knowledge Data Science During COVID-19 Pandemic](#)

[14.1 Introduction](#)

[14.2 Literature Review](#)

[14.3 Model Discussion](#)

[14.4 Results and Discussions](#)

[14.5 Conclusion](#)

[References](#)

[15 Semantic Data Science in the COVID-19 Pandemic](#)

[15.1 Crises Often Are Catalysts for New Technologies](#)

[15.2 The Domains of COVID-19 Semantic Data Science Research](#)

[15.3 Discussion](#)

[References](#)

[Index](#)

[Wiley End User License Agreement](#)

List of Tables

Chapter 1

[Table 1.1 Comparison of data science and business intelligence.](#)

Chapter 2

[Table 2.1 Common SQL commands.](#)

[Table 2.2 Common SQL operations.](#)

[Table 2.3 Common SQL order of execution.](#)

[Table 2.4 Worldwide electric energy consumption from 1981 to 2020.](#)

[Table 2.5 Sample dataset: prices of product \(X\) within 10 consecutive months.](#)

Chapter 3

[Table 3.1 Atoms and characteristics.](#)

Chapter 5

[Table 5.1 Data sets \[14\].](#)

[Table 5.2 Genes selected—AML/ALL data set \[14\].](#)

[Table 5.3 Genes selected - Lung Harvard2 data set \[14\].](#)

[Table 5.4 Signature genes related with various conditions for data set II \(Ovari...](#)

Chapter 6

[Table 6.1 Rotations, Ratios Max/Min, Execution times-Insert algorithm- Random ca...](#)

[Table 6.2 Rotations, Ratios Max/Min, Execution times-Insert algorithm- Random ca...](#)

[Table 6.3 Number of rotations, execution times-Delete algorithm-Random case.](#)

[Table 6.4 Insert algorithm—Ascending case.](#)

[Table 6.5 Execution time-Insert algorithm-Random case.](#)

[Table 6.6 Execution time-Delete algorithm-Random case.](#)

Chapter 7

[Table 7.1 Size of security ontologies.](#)

Chapter 9

[Table 9.1 Comparsion of our approah with other bench mark approach.](#)

Chapter 11

[Table 11.1 Richness of security ontologies via ontometric tool \(a\) security onto...](#)

Chapter 14

[Table 14.1 COVID-19 pandemic timeline from January 2020 to December 2020.](#)

[Table 14.2 COVID-19 pandemic timeline from January 2021 to June 2021.](#)

[Table 14.3 The country wise display of confirmed, recovered, deaths, and active ...](#)

[Table 14.4 14.4 \(a\), 14.4 \(b\), 14.4 \(c\), and 14.4 \(d\) depict the predicted outpu...](#)

[Table 14.5 Model performance for active cases.](#)

[Table 14.6 Model performance for confirmed cases.](#)

[Table 14.7 Model performance for recovered cases.](#)

[Table 14.8 Model performance for death cases.](#)

[Table 14.9 Display of top 20 countries confirmed and recovered cases as on 29.05...](#)

[Table 14.10 COVID-19 patients actual and predicted confirmed, deaths, recovered ...](#)

[Table 14.11 Forecast values for 5th day, 10th day, and 15th day new infected cas...](#)

Chapter 15

[Table 15.1 Summary table of all reviewed research.](#)

List of Illustrations

Chapter 1

[Figure 1.1 Proceeding of IJCAI-workshop.](#)

[Figure 1.2 Data science Venn diagram.](#)

[Figure 1.3 Job growth on analytics and data science.](#)

[Figure 1.4 Components of data science.](#)

[Figure 1.5 Features of SAS.](#)

[Figure 1.6 Features of Apache Spark.](#)

[Figure 1.7 Features of D3.js.](#)

[Figure 1.8 Features of MATLAB.](#)

[Figure 1.9 Features of excel.](#)

[Figure 1.10 Features of tableau.](#)

[Figure 1.11 Features of NLTK.](#)

[Figure 1.12 Features of TensorFlow.](#)

Chapter 2

[Figure 2.1 Diverse fields of data science.](#)

[Figure 2.2 Data scientist vs. data engineer \[5\].](#)

[Figure 2.3 Structured data vs. unstructured data vs semistructured data.](#)

[Figure 2.4 Top ten tools for data science.](#)

[Figure 2.5 Front view of Anaconda platform and Jupyter for IPython.](#)

[Figure 2.6 Companies that use R language for data analytics.](#)

[Figure 2.7 Operators of R language.](#)

[Figure 2.8 Most common operators of R language.](#)

[Figure 2.9 Example of schema for RDB of two tables.](#)

[Figure 2.10 Simplified example for outliers in the dataset.](#)

[Figure 2.11 Algorithm of outlier detection used in Tukey labeling.](#)

[Figure 2.12 Mouse pointer points to statistic chart symbol, and then points box ...](#)

[Figure 2.13 Box and whisker with single outlier.](#)

[Figure 2.14 Linear regression vs logistic regression.](#)

[Figure 2.15 Using scatter plot tool in microsoft excel.](#)

[Figure 2.16 The scatter plot for the energy data records provided in Table 2.4.](#)

[Figure 2.17 “Add Trend line” option for the scatter plot.](#)

[Figure 2.18 Using Linear regression via “Add Trend line” option for the scatter ...](#)

[Figure 2.19 Using forecasting tool in Microsoft Excel.](#)

[Figure 2.20 Preview of forecasting results.](#)

[Figure 2.21 Final forecasting results of six future values.](#)

[Figure 2.22 Apache spark architecture.](#)

[Figure 2.23 MongoDB data store architecture.](#)

[Figure 2.24 MATLAB computing system architecture.](#)

[Figure 2.25 Neo4j_graph database architecture.](#)

[Figure 2.26 VMWare virtualization system architecture.](#)

Chapter 3

[Figure 3.1 Periodic table of atoms. \(Sources: <https://en.wikipedia.org/wiki/Peri...>](#)

[Figure 3.2 Schema of data item, data field, record and file \(Sources: <https://ww...>](#)

Chapter 5

[Figure 5.1 The history of DNA sequencing technologies \(Source: <https://www.cd-ge...>](#)

[Figure 5.2 Block diagram of proposed method I.](#)

[Figure 5.3 Block diagram of proposed method II \[14\].](#)

Chapter 6

[Figure 6.1 P\(2\)-BST.](#)

[Figure 6.2 P\(3\)-BST.](#)

[Screenshot 6-01 AVL tree insert algorithm.](#)

[Screenshot 6-02 AVL tree insert algorithm.](#)

[Figure 6.3 Partitioning.](#)

[Figure 6.4 Departitioning.](#)

[Figure 6.5 Departitioning \(special case\).](#)

[Figure 6.6 Restructuring-partitioning.](#)

[Figure 6.7 Transforming.](#)

[Figure 6.8 Step-by-step insert process.](#)

[Figure 6.9 Step-by-step delete process.](#)

[Figure 6.10 P\(2\)-BST and P\(3\)-BST versus corresponding Red-Black trees.](#)

[Screenshot 6-03 PBST node implementation.](#)

[Screenshot 6-04 Red-Black tree node implementation.](#)

[Screenshot 6-05 AVL tree node implementation.](#)

[Figure 6.11 Numbers of rotations-Insert algorithm-Random case.](#)

[Figure 6.12 Execution times-Insert algorithm-Random case.](#)

[Figure 6.13 Number of rotations-Delete algorithm-Random case.](#)

[Figure 6.14 Execution times-Delete algorithm-Random case.](#)

[Figure 6.15 Maximal height of path-Insert algorithm-Ascending case.](#)

[Figure 6.16 Ratio Max/Min-Insert algorithm-Ascending case.](#)

[Figure 6.17 Execution time-Insert algorithm-Ascending case.](#)

[Figure 6.18 Distribution of class and simple nodes-Insert algorithm-Ascending ca...](#)

Chapter 7

[Figure 7.1 Languages leading to OWL.](#)

[Figure 7.2 From HTML to linked data.](#)

[Figure 7.3 Obtained articles per year \(from 2010 to July 2021\)_\(a\) Science direc...](#)

[Figure 7.4 Screening of the articles.](#)

[Figure 7.5 Classification framework for security ontologies analysis.](#)

[Figure 7.6 Evaluation of security ontology via OOPS! tool of InFra_OE framework.](#)

[Figure 7.7 Pitfall rate of security ontologies.](#)

Chapter 8

[Figure 8.1 Water used in coal mining.](#)

[Figure 8.2 Interfacing block representation of proposed architecture.](#)

[Figure 8.3 \(a\) Arduino Uno \(b\) Arduino Uno Atmega328p pin mapping.](#)

[Figure 8.4 LM 35 temperature sensor.](#)

[Figure 8.5 16X2 LCD display module.](#)

[Figure 8.6 \(a\) DHT11 temperature-humidity sensor. \(b\) DHT11 connection with proc...](#)

[Figure 8.7 Negative temperature coefficient of resistance.](#)

[Figure 8.8 Moisture sensor module.](#)

[Figure 8.9 MQ 135 gas sensor module.](#)

[Figure 8.10 SIM900 GSM GPRS Module connection with Arduino Uno.](#)

[Figure 8.11 Solar PV system.](#)

[Figure 8.12 LM-35 interfaces with Arduino.](#)

[Figure 8.13 DHT11 interface with Arduino.](#)

[Figure 8.14 MQ-X module interface with Arduino.](#)

[Figure 8.15 NEO 6M GPS module interface with Arduino.](#)

[Figure 8.16 DTH11 module pc display result.](#)

Chapter 9

[Figure 9.1 Inverted index stucture stucture with other indexed index structure.](#)

[Figure 9.2 Hierarchical structure of sentence classification.](#)

[Figure 9.3 Fuzzy system.](#)

Chapter 10

[Figure 10.1 \(a\), \(b\) and \(c\) Predict the cost of a house based on square feet.](#)

[Figure 10.2 Semantic stack architecture.](#)

[Figure 10.3 RDF graph.](#)

[Figure 10.4 Property values for Mary_Doe.](#)

[Figure 10.5 Explanation for inference that Mary_Doe has Aunt Sarah_Doe.](#)

Chapter 11

[Figure 11.1 Ontology evaluation tools.](#)

Chapter 13

[Figure 13.1 Calculation of the elements of the semantic matrix.](#)

[Figure 13.2 An example of calculating the strength of the semantic relationship ...](#)

[Figure 13.3 Relations of the distances \$r_i^k\(r_j^k\)\$ of objects i,j to the center of the cl...](#)

[Figure 13.4 Geometric representation of cluster k objects for the extended PDC c...](#)

Chapter 14

[Figure 14.1 Illustration of the SARS-CoV-2 virion.](#)

[Figure 14.2 Transmission stage of corona virus.](#)

[Figure 14.3 Flowchart of the implementation methodology.](#)

[Figure 14.4 World map of COVID-19 active cases as of 29.05.2021. The darker the ...](#)

[Figure 14.5 The graph depicts the observed date Vs data frame \(because of the hu...](#)

[Figure 14.6 \(a-b\) shows the actual and predicted cases active, death, recovered,...](#)

[Figure 14.6 \(c-d\) shows the actual and predicted cases active, death, recovered,...](#)

[Figure 14.7 Trend, weekly, and daily cases of active case.](#)

[Figure 14.8 Trend, weekly, and daily cases of confirmed cases.](#)

[Figure 14.9 Trend, weekly, and daily cases of recovered cases.](#)

[Figure 14.10 Trend, weekly, and daily cases of death cases.](#)

[Figure 14.11 Trend of the performance measures \(i\) MSE, \(ii\) RMSE, \(iii\) MAE for...](#)

[Figure 14.12 Trend of the performance measures \(i\) MSE, \(ii\) RMSE, \(iii\) MdAPE f...](#)

[Figure 14.13 Trend of the performance measures \(i\) MSE, \(ii\) RMSE, and \(iii\) MAP...](#)

[Figure 14.14 Trend of the performance measures \(i\) MSE, \(ii\) RMSE, and \(iii\) MAP...](#)

[Figure 14.15 The world's top 20 countries recovered vs confirmed cases as on 29....](#)

[Figure 14.16 Worldwide top countries daily COVID-19 cases of actual vs predicted...](#)

Chapter 15

[Figure 15.1 Semantic technology COVID-19 domains.](#)

[Figure 15.2 Amazon semantic search tool.](#)

[Figure 15.3 COVID*GRAPH semantic search tool.](#)

[Figure 15.4 Knowledge graph visualization of COVID-19 concepts and relevant pape...](#)

[Figure 15.5 Example of papers discovered via network graph semantic search.](#)

[Figure 15.6 The Johns Hopkins dashboard.](#)

[Figure 15.7 The NY times dataset in the Stardog cloud.](#)

[Figure 15.8 COVID trials visualization.](#)

[Figure 15.9 Drug repurposing visualization.](#)

[Figure 15.10 Infection paths visualized by CODO.](#)

[Figure 15.11 Geographic information visualized in CODO.](#)

Scrivener Publishing

100 Cummings Center, Suite 541J

Beverly, MA 01915-6106

Publishers at Scrivener

Martin Scrivener (martin@scrivenerpublishing.com)

Phillip Carmical (pcarmical@scrivenerpublishing.com)

Data Science with Semantic Technologies

Theory, Practice, and Application

Edited by

Archana Patel

*Department of Software Engineering, School of Computing
and Information Technology, Eastern International
University, Vietnam*

Narayan C. Debnath

*School of Computing and Department of Computer Science
and Engineering, School of Engineering Vietnam*

and

Bharat Bhusan

*Technology, Sharda University Information Technology,
India*



WILEY

This edition first published 2022 by John Wiley & Sons, Inc., 111 River Street, Hoboken, NJ 07030, USA and Scrivener Publishing LLC, 100 Cummings Center, Suite 541J, Beverly, MA 01915, USA © 2022 Scrivener Publishing LLC

For more information about Scrivener publications please visit

www.scrivenerpublishing.com.

All rights reserved. No part of this publication may be reproduced, stored in a retrieval system, or transmitted, in any form or by any means, electronic, mechanical, photocopying, recording, or otherwise, except as permitted by law. Advice on how to obtain permission to reuse material from this title is available at <http://www.wiley.com/go/permissions>.

Wiley Global Headquarters

111 River Street, Hoboken, NJ 07030, USA

For details of our global editorial offices, customer services, and more information about Wiley products visit us at www.wiley.com.

Limit of Liability/Disclaimer of Warranty

While the publisher and authors have used their best efforts in preparing this work, they make no representations or warranties with respect to the accuracy or completeness of the contents of this work and specifically disclaim all warranties, including without limitation any implied warranties of merchantability or fitness for a particular purpose. No warranty may be created or extended by sales representatives, written sales materials, or promotional statements for this work. The fact that an organization, website, or product is referred to in this work as a citation and/or potential source of further information does not mean that the publisher and authors endorse the information or services the organization, website, or product may provide or recommendations it may make. This work is sold with the understanding that the publisher is not engaged in rendering professional services. The advice and strategies contained herein may not be suitable for your situation. You should consult with a specialist where appropriate. Neither the publisher nor authors shall be liable for any loss of profit or any other commercial damages, including but not limited to special, incidental, consequential, or other damages. Further, readers should be aware that websites listed in this work may have changed or disappeared between when this work was written and when it is read.

Library of Congress Cataloging-in-Publication Data

ISBN 9781119864981

Cover image: PixaBay.Com

Cover design by Russell Richardson

Set in size of 11pt and Minion Pro by Manila Typesetting Company, Makati, Philippines

Printed in the USA

10 9 8 7 6 5 4 3 2 1

Preface

Data Science is an invaluable resource that deals with vast volumes of data using modern tools and techniques to find unseen patterns, derive meaningful information, and make business decisions. To create intelligence in data science, it becomes necessary to utilize the semantic technologies which allow machine-readable representation of data. This intelligence uniquely identifies and connects data with common business terms, and also enables users to communicate with data. Instead of structuring the data, semantic technologies help users to understand the meaning of the data by using the concepts of semantics, ontology, OWL, linked data, and knowledge graphs. These technologies assist organizations in understanding all of the stored data, adding value to it, and enabling insights that were not available before. Organizations are also using semantic technologies to unearth precious nuggets of information from vast volumes of data and to enable more flexible use of data. These technologies can deal with the existing problems of data scientists and help them in making better decisions for any organization. All of these needs are part of a focused shift towards utilization of semantic technologies in data science that provide knowledge along with the ability to understand, reason, plan, and learn with existing and new data sets. These technologies also generate expected, reproducible, user-desired results.

This book aims to provide a roadmap for the deployment of semantic technologies in the field of data science. Moreover, it highlights how data science enables the user to create intelligence through these technologies by exploring the opportunities and eradicating the challenges

in current and future time frames. It can serve as an important guide to applications of data science with semantic technologies for the upcoming generation and thus is a unique resource for scholars, researchers, professionals and practitioners in this field. Following is a brief description of the subjects covered in the 15 chapters of the book.

- [Chapter 1](#) provides a brief introduction to data science. It addresses various aspects of data science such as what a data scientist does and why data science is in demand; the history of data science and how it differs from business intelligence; the life cycle of data science and data science components; why data science is important; the challenges of data science; the tools used for data science; and the benefits and applications of data science.
- [Chapter 2](#) provides an overview of the top 10 tools and applications that should be of interest to any data scientist. Its objective includes, but is not limited to, realizing the use of Python in developing solutions to data science tasks; recognizing the use of R language as an open-source data science provider; traveling around the SQL to provide structured models for data science projects; navigating through data analytics and statistics using Excel; and using D3.js scripting tools for data visualization. Also, practical examples/case studies are provided on data visualization, data analytics, regression, forecasting, and outlier detection.
- [Chapter 3](#) presents the use of data modeling for data science, revealing the possibility of a new side of the data. The chapter covers different types of data (unstructured data, semi-structured data, structured data, hybrid (un/semi)-structured data and big data) and data model design.

- [Chapter 4](#) shows data management by considering language based on the novelty view of data. The chapter focuses on data life cycle, data distribution and CAP theorem.
- [Chapter 5](#) presents the role of data science in healthcare. There are several fields in the healthcare sector, such as predictive modeling, genetics, etc., which make use of data science for diagnosis and drug discovery, thereby increasing usability of precision medicine.
- [Chapter 6](#) provides a new balanced binary search tree that generates two kinds of nodes: simple and class nodes. Two advantages make the new structure attractive. First, it subsumes the most popular data structures of AVL and Red-Black trees. Second, it proposes other unknown balanced binary search trees in which we can adjust the maximal height of paths between $1.44 \lg(n)$ and $2 \lg(n)$, where n is the number of nodes in the tree and \lg the base-two logarithm.
- [Chapter 7](#) shows the study of machine learning and deep learning algorithms with detailed and analytical comparisons, which help new and inexperienced medical professionals or researchers in the medical field. The proposed machine learning model has an accurate algorithm that works with rich healthcare data, a high-dimensional data handling system, and an intelligent framework that uses different data sources to predict heart disease. This chapter uses an ensemble-based deep learning model with optimal feature selection to improve accuracy.
- [Chapter 8](#) presents an IoT-based automated fire control system in a mining area which will help to protect many valuable lives whenever an accident occurs due to fire. In the experimental application, different types of sensors for

temperature, moisture, and gas are used to sense the different environmental data.

- [Chapter 9](#) offers an aspect identification method for sentiment sentences in review documents. The chapter describes two key tasks—one for extracting significant features from the reviews and another for identification of degrees of product reviews.
- [Chapter 10](#) shows the research that paved the way for semantic technology. It then describes each of the semantic pillars with examples and explanations of the business value of each technology.
- [Chapter 11](#) describes the ontology evaluation tools and then focuses on the evaluation of the security ontologies. The existing ontology evaluation tools are classified under two categories; namely, domain-dependent ontology evaluation tools and domain-independent ontology evaluation tools. The evaluation of security ontology assesses the quality of ontology among the available ontologies.
- [Chapter 12](#) discusses the main concepts of health data, data science, health data science, examples of the application of health data science and related challenges. In addition, it also highlights the application of semantic technologies in health data science and the challenges that lie ahead of using these technologies.
- [Chapter 13](#) proposes an original hybrid optimization approach based on two different mixed integer programming statements. The first statement is based on minimizing the sum of pairwise distances between all objects (PDC clustering), while the second statement is based on minimizing the total distance from objects to cluster centers (CC clustering). Computational experiments showed that the hybrid method developed for solving the

clustering problem combines the advantages of both approaches—the speed of the k-means method and the accuracy of PDC clustering—which makes it possible to get rid of the main drawback of the k-means, namely, the lack of guaranteed determining of the global optimum.

- [Chapter 14](#) uses a model for the analysis of time series data which highly depend on the novel coronavirus 2019. This model predicts the future trend of confirmed, recovered, active, and death cases based on the available data from January 22, 2020 to May 29, 2021. The present model predicted the spread of COVID-19 for a future period of 30 days. The RMSE, MSE, MAE, and MdAPE metrics are used for the model evaluation.

- [Chapter 15](#) focuses on systems that incorporated real-world data utilized by actual users. It first describes a new methodology for the survey and then covers the various domains where semantic technology can be applied and some of the most impressive systems developed in each domain.

Finally, the editors would like to sincerely thank all the authors and reviewers who contributed their time and expertise for the successful completion of this book. The editors also appreciate the support given by Scrivener Publishing, which allowed us to compile the first edition of this book.

The Editors
Archana Patel
Narayan C. Debnath
Bharat Bhusan
June 2022

1

A Brief Introduction and Importance of Data Science

Karthika N.^{[1*](#)}, Sheela J.^{[1](#)} and Janet B.^{[2](#)}

^{[1](#)}*Department of SCOPE, VIT-AP University, Amaravati, Andhra Pradesh, India*

^{[2](#)}*Department of Computer Applications, National Institute of Technology, Tiruchirappalli, India*

Abstract

Data is very important component of any organization. According to International Data Corporation, by 2025, global data will reach to 175 zettabytes. They need data to help them make careful decisions in business. Data is worthless until it is transformed into valuable data. Data science plays a vital role in processing and interpreting data. It focuses on the analysis and management of data too. It is concerned with obtaining useful information from large datasets. It is frequently applied in a wide range of industries, including healthcare, marketing, banking, finance, policy work, and more. This enables companies to make informed decisions around growth, optimization, and performance. In this brief monograph, we address following questions.

What is data science and what does a data scientist do? Why data science is in demand? History of data science, how data science differs from business intelligence? The lifecycle of data science, data science components, why data science is important? Challenges of data science, tools used for data science, benefits and applications of data science.

Keywords: Data science, history, lifecycle, components, tools

1.1 What is Data Science? What Does a Data Scientist Do?

Data is very important component of any organization. According to International Data Corporation, by 2025, global data will reach to 175 zettabytes. They need data to help them make careful decisions in business. Data is worthless until it is transformed into valuable data. Data science plays a vital role in processing and interpreting

data. It focuses on the analysis and management of data too. It is concerned with obtaining useful information from large datasets. It is frequently applied in a wide range of industries, including healthcare, marketing, banking, finance, policy work, and more. This enables companies to make informed decisions around growth, optimization, and performance. In nutshell, Data science is an integrative strategy for deriving actionable insights from today's organizations' massive and ever-increasing data sets. Preparing data for analysis and processing, performing advanced data analysis, and presenting the findings to expose trends and allow stakeholders to make educated decisions are all part of data science [1, 2]. Data science experts are both well-known, data-driven individuals with advanced technical capabilities who can construct complicated quantitative algorithms to organize and interpret huge amounts of data in order to address questions and drive strategy in their company. This is combined with the communication and leadership skills required to provide tangible results to numerous stakeholders throughout a company or organization. Data scientists must be inquisitive and results-driven, with great industry-specific expertise and communication abilities that enable them to convey highly technical outcomes to non-technical colleagues. To create and analyze algorithms, they have a solid quantitative background in statistics and linear algebra, as well as programming experience with a focus on data warehousing, mining, and modeling [3].

1.2 Why Data Science is in Demand?

Data science is the branch of science concerned with the discovery, analysis, modeling, and extraction of useful information which has become a buzz in a lot of companies. Firms are increasingly aware that they have been sitting on

data treasure mines the priority with which this data must be analyzed, and ROI generated is obvious. We look at the most important reasons that data science professions are in high demand [4].

- **Data Organization**

During the mid-2000s IT boom, the emphasis was on “lifting and shifting” offline business operations into automated computer systems. Digital content generation, transactional data processing, and data log streams have all been consistent throughout the last two decades. This indicates that every company now has a plethora of information that it believes can really be valuable but does not know how to use. This is apparent in Glassdoor’s recent analysis, which identifies the 50 greatest jobs in modern era.

- **Scarcity of Trained Manpower**

According to a McKinsey Global Institute study, by 2018, the United States will be short 190,000 data scientists, 1.5 million managers, including analysts who would properly comprehend and make judgments based on Big Data. The need is particularly great in India, where the tools and techniques are available but there are not enough qualified people. Data scientists, who can perform analytics, and analytics consultants, who can analyze and apply data, are two sorts of talent shortages, according to Srikanth Velamakanni, co-founder and CEO of Fractal Analytics. The supply of talent in these fields, particularly data scientists, is extremely limited, and the demand is enormous.”

- **The Pay Is Outstanding**

A data science position is currently one of the highest paying in the market. The national average income for a data scientist/analyst in the United States, according