

DATA SCIENCE WITH SEMANTIC TECHNOLOGIES

THEORY, PRACTICE AND APPLICATION

Edited by

Archana Patel
Narayan C. Debnath
Bharat Bhusan

 Scrivener
Publishing

WILEY

Data Science with Semantic Technologies

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

Contents

Preface	xv
1 A Brief Introduction and Importance of Data Science	1
<i>Karthika N., Sheela J. and Janet B.</i>	
1.1 What is Data Science? What Does a Data Scientist Do?	2
1.2 Why Data Science is in Demand?	2
1.3 History of Data Science	4
1.4 How Does Data Science Differ from Business Intelligence?	9
1.5 Data Science Life Cycle	11
1.6 Data Science Components	13
1.7 Why Data Science is Important	14
1.8 Current Challenges	15
1.8.1 Coordination, Collaboration, and Communication	16
1.8.2 Building Data Analytics Teams	16
1.8.3 Stakeholders vs Analytics	17
1.8.4 Driving with Data	17
1.9 Tools Used for Data Science	19
1.10 Benefits and Applications of Data Science	28
1.11 Conclusion	28
References	29
2 Exploration of Tools for Data Science	31
<i>Qasem Abu Al-Haija</i>	
2.1 Introduction	32
2.2 Top Ten Tools for Data Science	35
2.3 Python for Data Science	35
2.3.1 Python Datatypes	36
2.3.2 Helpful Rules for Python Programming	37
2.3.3 Jupyter Notebook for IPython	37
2.3.4 Your First Python Program	38
2.4 R Language for Data Science	39

2.4.1	R Datatypes	39
2.4.2	Your First R Program	41
2.5	SQL for Data Science	44
2.6	Microsoft Excel for Data Science	48
2.6.1	Detection of Outliers in Data Sets Using Microsoft Excel	48
2.6.2	Regression Analysis in Excel Using Microsoft Excel	50
2.7	D3.JS for Data Science	57
2.8	Other Important Tools for Data Science	58
2.8.1	Apache Spark Ecosystem	58
2.8.2	MongoDB Data Store System	60
2.8.3	MATLAB Computing System	62
2.8.4	Neo4j for Graphical Database	63
2.8.5	VMWare Platform for Virtualization	65
2.9	Conclusion	66
	References	68
3	Data Modeling as Emerging Problems of Data Science	71
	<i>Mahyuddin K. M. Nasution and Marischa Elveny</i>	
3.1	Introduction	72
3.2	Data	72
3.2.1	Unstructured Data	74
3.2.2	Semistructured Data	74
3.2.3	Structured Data	76
3.2.4	Hybrid (Un/Semi)-Structured Data	77
3.2.5	Big Data	78
3.3	Data Model Design	79
3.4	Data Modeling	81
3.4.1	Records-Based Data Model	81
3.4.2	Non-Record-Based Data Model	84
3.5	Polyglot Persistence Environment	87
	References	88
4	Data Management as Emerging Problems of Data Science	91
	<i>Mahyuddin K. M. Nasution and Rahmad Syah</i>	
4.1	Introduction	92
4.2	Perspective and Context	92
4.2.1	Life Cycle	93
4.2.2	Use	95
4.3	Data Distribution	98
4.4	CAP Theorem	100

4.5	Polyglot Persistence	101
	References	102
5	Role of Data Science in Healthcare	105
	<i>Anidha Arulanantham, A. Suresh and Senthil Kumar R.</i>	
5.1	Predictive Modeling—Disease Diagnosis and Prognosis	106
5.1.1	Supervised Machine Learning Models	107
5.1.2	Clustering Models	110
5.1.2.1	Centroid-Based Clustering Models	110
5.1.2.2	Expectation Maximization (EM) Algorithm	110
5.1.2.3	DBSCAN	111
5.1.3	Feature Engineering	111
5.2	Preventive Medicine—Genetics/Molecular Sequencing	111
5.2.1	Technologies for Sequencing	113
5.2.2	Sequence Data Analysis with BioPython	114
5.2.2.1	Sequence Data Formats	114
5.2.2.2	BioPython	117
5.3	Personalized Medicine	121
5.4	Signature Biomarkers Discovery from High Throughput Data	122
5.4.1	Methodology I — Novel Feature Selection Method with Improved Mutual Information and Fisher Score	123
5.4.1.1	Algorithm for the Novel Feature Selection Method with Improved Mutual Information and Fisher Score	124
5.4.1.2	Computing F-Score Values for the Features	125
5.4.1.3	Block Diagram for the Method-1	125
5.4.1.4	Data Set	126
5.4.1.5	Identification of Biomarkers Using the Feature Selection Technique-I	127
5.4.2	Feature Selection Methodology-II — Entropy Based Mean Score with mRMR	128
5.4.2.1	Algorithm for the Feature Selection Methodology-II	130
5.4.2.2	Introduction to mRMR Feature Selection	132
5.4.2.3	Data Sets	132
5.4.2.4	Identification of Biomarkers Using Rank Product	133
5.4.2.5	Fold Change Values	133

Conclusion	136
References	136
6 Partitioned Binary Search Trees (P(h)-BST): A Data Structure for Computer RAM	139
<i>Pr. D.E Zegour</i>	
6.1 Introduction	140
6.2 P(h)-BST Structure	141
6.2.1 Preliminary Analysis	143
6.2.2 Terminology and Conventions	143
6.3 Maintenance Operations	143
6.3.1 Operations Inside a Class	145
6.3.2 Operations Between Classes (Outside a Class)	148
6.4 Insert and Delete Algorithms	153
6.4.1 Inserting a New Element	153
6.4.2 Deleting an Existing Element	157
6.5 P(h)-BST as a Generator of Balanced Binary Search Trees	160
6.6 Simulation Results	162
6.6.1 Data Structures and Abstract Data Types	164
6.6.2 Analyzing the Insert and Delete Process in Random Case	164
6.6.3 Analyzing the Insert Process in Ascending (Descending) Case	168
6.6.4 Comparing P(2)-BST/P(∞)-BST to Red-Black/AVL Trees	174
6.7 Conclusion	175
Acknowledgments	176
References	176
7 Security Ontologies: An Investigation of Pitfall Rate	179
<i>Archana Patel and Narayan C. Debnath</i>	
7.1 Introduction	179
7.2 Secure Data Management in the Semantic Web	184
7.3 Security Ontologies in a Nutshell	187
7.4 InFra_OE Framework	189
7.5 Conclusion	193
References	193
8 IoT-Based Fully-Automated Fire Control System	199
<i>Lalit Mohan Satapathy</i>	
8.1 Introduction	200
8.2 Related Works	201

8.3	Proposed Architecture	203
8.4	Major Components	205
8.4.1	Arduino UNO	205
8.4.2	Temperature Sensor	207
8.4.3	LCD Display (16X2)	208
8.4.4	Temperature Humidity Sensor (DHT11)	209
8.4.5	Moisture Sensor	210
8.4.6	CO ₂ Sensor	211
8.4.7	Nitric Oxide Sensor	212
8.4.8	CO Sensor (MQ-9)	212
8.4.9	Global Positioning System (GPS)	212
8.4.10	GSM Modem	213
8.4.11	Photovoltaic System	214
8.5	Hardware Interfacing	216
8.6	Software Implementation	218
8.7	Conclusion	222
	References	223
9	Phrase Level-Based Sentiment Analysis Using Paired Inverted Index and Fuzzy Rule	225
	<i>Sheela J., Karthika N. and Janet B.</i>	
9.1	Introduction	226
9.2	Literature Survey	228
9.3	Methodology	233
9.3.1	Construction of Inverted Wordpair Index	234
9.3.1.1	Sentiment Analysis Design Framework	235
9.3.1.2	Sentiment Classification	236
9.3.1.3	Preprocessing of Data	237
9.3.1.4	Algorithm to Find the Score	240
9.3.1.5	Fuzzy System	240
9.3.1.6	Lexicon-Based Sentiment Analysis	241
9.3.1.7	Defuzzification	242
9.3.2	Performance Metrics	243
9.4	Conclusion	244
	References	244
10	Semantic Technology Pillars: The Story So Far	247
	<i>Michael DeBellis, Jans Aasman and Archana Patel</i>	
10.1	The Road that Brought Us Here	248
10.2	What is a Semantic Pillar?	249
10.2.1	Machine Learning	249
10.2.2	The Semantic Approach	250

10.3	The Foundation Semantic Pillars: IRI's, RDF, and RDFS	252
10.3.1	Internationalized Resource Identifier (IRI)	254
10.3.2	Resource Description Framework (RDF)	254
10.3.2.1	Alternative Technologies to RDF: Property Graphs	256
10.3.3	RDF Schema (RDFS)	257
10.4	The Semantic Upper Pillars: OWL, SWRL, SPARQL, and SHACL	259
10.4.1	The Web Ontology Language (OWL)	260
10.4.1.1	Axioms to Define Classes	262
10.4.1.2	The Open World Assumption	263
10.4.1.3	No Unique Names Assumption	263
10.4.1.4	Serialization	264
10.4.2	The Semantic Web Rule Language	264
10.4.2.1	The Limitations of Monotonic Reasoning	267
10.4.2.2	Alternatives to SWRL	267
10.4.3	SPARQL	268
10.4.3.1	The SERVICE Keyword and Linked Data	268
10.4.4	SHACL	271
10.4.4.1	The Fundamentals of SHACL	272
10.5	Conclusion	274
	References	274
11	Evaluating Richness of Security Ontologies for Semantic Web	277
	<i>Ambrish Kumar Mishra, Narayan C. Debnath and Archana Patel</i>	
11.1	Introduction	277
11.2	Ontology Evaluation: State-of-the-Art	280
11.2.1	Domain-Dependent Ontology Evaluation Tools	281
11.2.2	Domain-Independent Ontology Evaluation Tools	282
11.3	Security Ontology	284
11.4	Richness of Security Ontologies	287
11.5	Conclusion	295
	References	295
12	Health Data Science and Semantic Technologies	299
	<i>Haleh Ayatollahi</i>	
12.1	Health Data	300
12.2	Data Science	301
12.3	Health Data Science	301
12.4	Examples of Health Data Science Applications	304
12.5	Health Data Science Challenges	306

12.6	Health Data Science and Semantic Technologies	308
12.6.1	Natural Language Processing (NLP)	309
12.6.2	Clinical Data Sharing and Data Integration	310
12.6.3	Ontology Engineering and Quality Assurance (QA)	311
12.7	Application of Data Science for COVID-19	313
12.8	Data Challenges During COVID-19 Outbreak	314
12.9	Biomedical Data Science	315
12.10	Conclusion	316
	References	317
13	Hybrid Mixed Integer Optimization Method for Document Clustering Based on Semantic Data Matrix	323
	<i>Tatiana Avdeenko and Yuri Mezentsev</i>	
13.1	Introduction	324
13.2	A Method for Constructing a Semantic Matrix of Relations Between Documents and Taxonomy Concepts	327
13.3	Mathematical Statements for Clustering Problem	330
13.3.1	Mathematical Statements for PDC Clustering Problem	330
13.3.2	Mathematical Statements for CC Clustering Problem	334
13.3.3	Relations between PDC Clustering and CC Clustering	336
13.4	Heuristic Hybrid Clustering Algorithm	340
13.5	Application of a Hybrid Optimization Algorithm for Document Clustering	342
13.6	Conclusion	344
	Acknowledgment	344
	References	344
14	Role of Knowledge Data Science During COVID-19 Pandemic	347
	<i>Veena Kumari H. M. and D. S. Suresh</i>	
14.1	Introduction	348
14.1.1	Global Health Emergency	350
14.1.2	Timeline of the COVID-19	351
14.2	Literature Review	354
14.3	Model Discussion	356
14.3.1	COVID-19 Time Series Dataset	357
14.3.2	FBProphet Forecasting Model	358
14.3.3	Data Preprocessing	360

14.3.4	Data Visualization	360
14.4	Results and Discussions	362
14.4.1	Analysis and Forecasting: The World	362
14.4.2	Performance Metrics	371
14.4.3	Analysis and Forecasting: The Top 20 Countries	377
14.5	Conclusion	388
	References	389
15	Semantic Data Science in the COVID-19 Pandemic	393
	<i>Michael DeBellis and Biswanath Dutta</i>	
15.1	Crises Often Are Catalysts for New Technologies	393
15.1.1	Definitions	394
15.1.2	Methodology	395
15.2	The Domains of COVID-19 Semantic	
	Data Science Research	397
15.2.1	Surveys	398
15.2.2	Semantic Search	399
15.2.2.1	Enhancing the CORD-19 Dataset with Semantic Data	399
15.2.2.2	CORD-19-on-FHIR -- Semantics for COVID-19 Discovery	400
15.2.2.3	Semantic Search on Amazon Web Services (AWS)	400
15.2.2.4	COVID*GRAPH	402
15.2.2.5	Network Graph Visualization of CORD-19	403
15.2.2.6	COVID-19 on the Web	404
15.2.3	Statistics	405
15.2.3.1	The Johns Hopkins COVID-19 Dashboard	405
15.2.3.2	The NY Times Dataset	406
15.2.4	Surveillance	406
15.2.4.1	An IoT Framework for Remote Patient Monitoring	406
15.2.4.2	Risk Factor Discovery	408
15.2.4.3	COVID-19 Surveillance in a Primary Care Network	408
15.2.5	Clinical Trials	409
15.2.6	Drug Repurposing	411
15.2.7	Vocabularies	414
15.2.8	Data Analysis	415

15.2.8.1	CODO	415
15.2.8.2	COVID-19 Phenotypes	416
15.2.8.3	Detection of “Fake News”	417
15.2.8.4	Ontology-Driven Weak Supervision for Clinical Entity Classification	417
15.2.9	Harmonization	418
15.3	Discussion	418
15.3.1	Privacy Issues	420
15.3.2	Domains that May Currently be Under Utilized	421
15.3.2.1	Detection of Fake News	421
15.3.2.2	Harmonization	421
15.3.3	Machine Learning and Semantic Technology: Synergy Not Competition	422
15.3.4	Conclusion	423
	Acknowledgment	423
	References	423
	Index	427

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

A Brief Introduction and Importance of Data Science

Karthika N.^{1*}, Sheela J.¹ and Janet B.²

¹Department of SCOPE, VIT-AP University, Amaravati, Andhra Pradesh, India

²Department 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

*Corresponding author: bharathikarthika@gmail.com

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 to Glass Door, is more than \$62,000. In India, pay is heavily influenced by experience. Those with the appropriate skillset can earn up to 19 LPA. (source: PayScale.)

- **The “X” Factor**

A data scientist’s major responsibility are exceptional and specific to the position. Because of nature of the profession, they may flourish in their careers by integrating several analytical expertise across diverse areas such as big data, machine learning, and so on. This vast knowledge base gives them an unsurpassed reputation or X-factor.

- **Data Scientists’ Democratization**

Tech behemoths are not the only ones who need data scientists. According to a Harvard Business Report issued many years ago, “Organizations in the top list of their area in the

use of data-driven decision making were, on average, 5% more productive and 6% more profitable than their peers”. Even mid-sized and small organizations have been driven to adopt data science because of this. In truth, many small businesses are trying to hire entry-level data scientists for a fair wage. This works well for both. The scientist will be able to further develop his or her skills, and the company will be able to pay less than it would otherwise.

- **Fewer Barriers for Professionals**

Data science is open to a wide range of experts from varied backgrounds because it is a relatively new discipline. Math/statistics, computer science and engineering, and natural science are all areas of knowledge for today’s data scientists. Some perhaps have social science, economics, or business degrees. They have all devised a problem-solving technique and improved their skills through formal or online education.

- **Abundance of Jobs**

Data science is employed in a wide range of business sectors, from production to healthcare, Information Technology to finance, therefore there are plenty of data science jobs available for individuals who are interested and willing to put in the effort. It is true not only in terms of industries, but also in terms of geography. So, regardless of one’s geographical location or current domain, data science and analytics are available to everybody.

- **A Wide Range of Roles**

Even if data science job is indeed a broad term, there are numerous subroles that fall under its scope. Data scientists, data architects, business intelligence engineers, business analysts, data engineers, database administrators, and data analytics managers are all in considerable demand.

1.3 History of Data Science

The terminology “data science” was just recently coined a new profession interested in trying to make sense of large volumes of data. Making sense of data, on the other hand, has a significant background, and it has been addressed for years by many computer scientists, scientists, librarians, statisticians, and others. The history below shows how the terminology:

data science” evolved over time, as well as attempts to describe it and associated concepts [5].

In 1974, Peter Naur’s book gives a broad overview of modern data processing techniques that are employed in a variety of applications. The IFIP Guide to Data Processing Concepts and Terms states that it is organized around the data principle: “Data is a codified representation of ideas or facts that may be communicated and even perhaps changed by certain process.” According to the book’s preface, a course plan titled “Datalogy, the science of data and data processes, and its position in education” was presented at the 1968 IFIP Congress, and the name “data science” has been widely used since then. Data science, according to Naur, is defined as “the science of working with data after it has been established, but the relationship of the data to what it represents is assigned to other disciplines and sciences.”

In 1977, the International Association for Statistical Computing (IASC) was founded as an ISI chapter. “The goal of the IASC is to connect conventional statistical techniques, innovative computer technology, and domain specialists’ skills to transform data into knowledge and information,” says the organization.” 1989 The first Knowledge Discovery in Databases (KDD) workshop is arranged and chaired by Gregory Piatetsky-Shapiro. Figure 1.1 shows the proceeding of IJCAI-Workshop. It was renamed the ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD) in 1995. September 1994, an article on “Database Marketing” appears in *Business Week*: “Enterprises are acquiring vast amounts of data on you, processing it to determine how likely someone really is to purchase a product, after then using that intelligence to design a marketing strategy perfectly tuned to find a way to convince you to do so... A prior spike of anticipation in the 1980s, sparked by the extensive use of checkout scanners, resulted in severe disappointment. Several organizations were overwhelmed with vast amount of data and were unable to do anything valuable with it... Despite this, many corporations recognize that they have no other option except to enter the database marketing arena.”

In 1997, the journal *Data Mining and Knowledge Discovery* is set up, with the reversal of the two terms in the title emphasizing the rise of “data mining” as the standard term for “extracting information from vast data-sets.” December 1999, “Existing statistical procedures perform effectively with relatively small data sets, Jacob Zahavi says in *Knowledge @ Wharton*’s “Mining Data for Nuggets of Knowledge.” Today’s databases, on the other hand, can have trillions of rows and columns of data. In data mining, scalability is a major concern. Another technical issue is creating models that would better analyze data, recognize nonlinear relations, and interaction

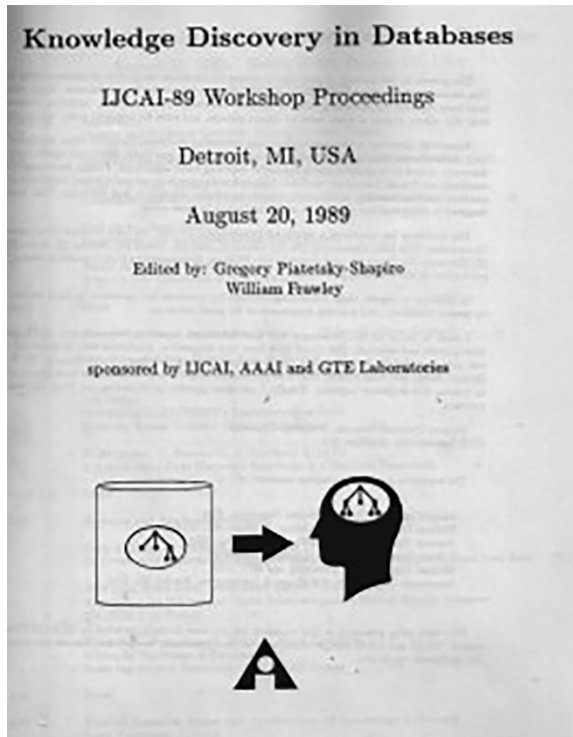


Figure 1.1 Proceeding of IJCAI-workshop.

among elements. To handle web-site issues, specialized data mining techniques may need to be developed.”

2001 “Data Science: An Action Plan for Expanding the Technical Areas of the Field of Statistics,” by William S. Cleveland, is published. That is a plan to “expand the key areas of technological endeavor in statistics.” The new area will be dubbed “data science” because the notion is ambitious and requires significant change. Cleveland compares and contrasts the proposed new discipline to computer science and current data mining research. 2001 “Statistical Modeling: The Two Cultures,” by Leo Breiman, is published. When it comes to utilizing statistical modeling to draw conclusions from data, there are two contrasting cultures. A stochastic data model is assumed to have generated the data.

The statistical community has been devoted to using data models nearly exclusively. This emphasis has resulted in useless theory, spurious findings, and the exclusion of statisticians from a wide spectrum of current situations. Algorithmic modeling has advanced significantly in theory and in

practice in domains other than statistics that may be used on huge challenging as well as on small datasets as the more useful and consistent alternative to data modeling. If we want to utilize data to resolve problems as a domain, we ought to move away from merely using data models and use a wider range of tools. The International Council for Science's Committee on Data for Science and Technology (CODATA) publishes the journal (ICSU). January in the year 2003 "By data science, people mean practically the whole thing that has to do with data: acquiring, examining, and modeling," according to the launch of the Journal of Data Science. However, the most important feature is its applicability—it may be used in a wide range of situations. The journal is primarily concerned with the application of statistical methods in general. All data scientists will be able to submit their viewpoints and ideas in the Journal of Data Science. The data science venn diagram is depicted by Figure 1.2.

In "A Taxonomy of Data Science," Chris Wiggins and Hilary Mason write in September 2010: One potential taxonomy... of what a data scientist does may be: obtain, scrub, explore, model, and perceive, in approximately chronological order. Data science is evidently a merger of the hacker arts of statistics and machine learning, as well as math and data topic knowledge for the analysis to be understandable. It necessitates creative thinking and a willingness to learn. "To become a properly trained data scientist, one needs comprehend a great deal," writes Drew Conway in "The Data Science

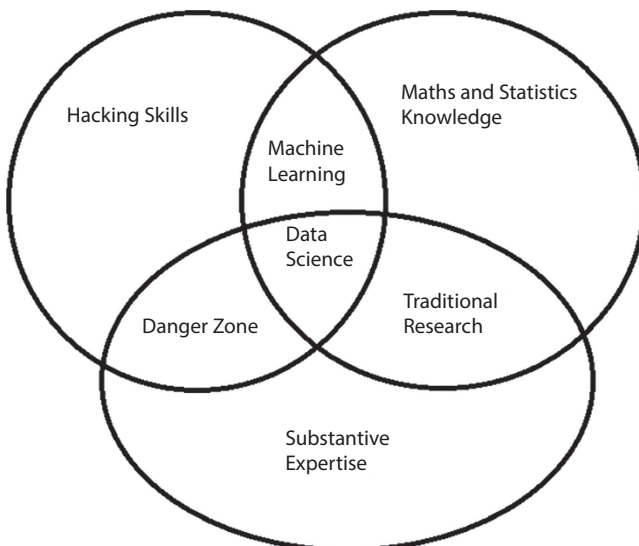


Figure 1.2 Data science Venn diagram.

Venn Diagram” from September 2010. Regrettably, simply citing literature and teachings is insufficient to untangle the tangles.

In the month of May 2011, in his article “Why the phrase ‘data science’ is wrong but useful,” Pete Warden writes: According to P., there is no commonly agreed-upon boundary between what is inside and beyond the purview of data science. Rather than choosing a topic in the beginning and later gathering data to shed light on it, they tend to be more concerned with what the data can disclose and then picking fascinating strands to pursue. Matthew J. Graham presented “The Art of Data Science” at the Astro statistics and Data Mining in Large Astronomical Databases workshop in June 2011. “We need to learn new abilities to succeed in the modern data-intensive world of twenty-first-century science,” he says. “We

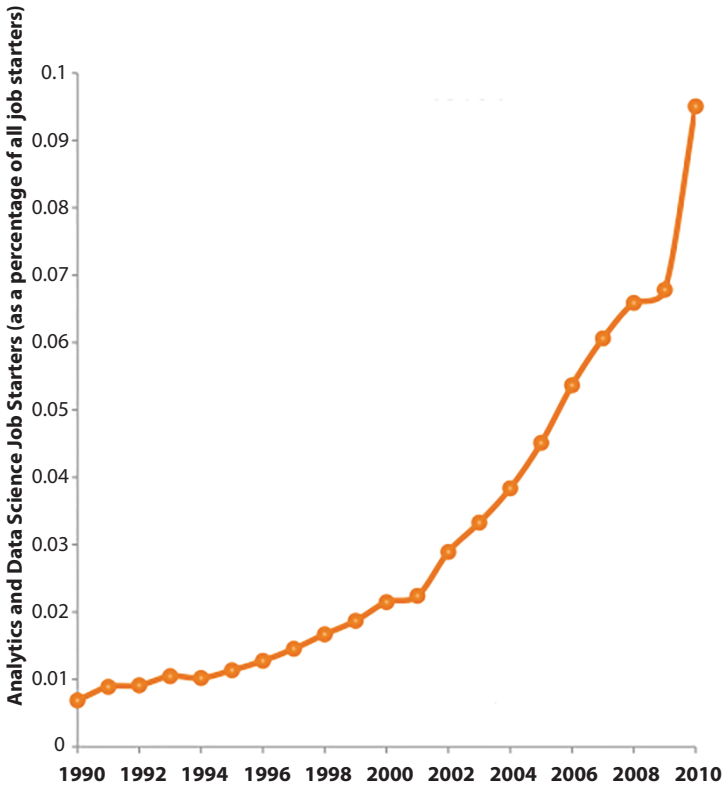


Figure 1.3 Job growth on analytics and data science.

should understand how [data] is perceived and expressed, as well as how it links to physical space and time.”

The title of “business analyst” felt too constrictive. The title of “data analyst” was a candidate, but we believed it limited what people could do. After all, many of the members of our teams were highly skilled engineers. The phrase “research scientist” was used by firms, including Sun, Yahoo, HP, IBM, and Xerox as a fitting job title. We assumed, however, that the majority of research scientists worked on future and abstract issues in labs apart from the product development groups. It would take ages for lab research to have an impact on big items. As a replacement, our teams concentrated on developing data applications that can have an instant and significant influence on the firm. Data scientist seems to be the greatest fit: individuals who combine data and research to produce something new [6]. Figure 1.3 shows the job growth on analytics and data science.

1.4 How Does Data Science Differ from Business Intelligence?

Business intelligence (BI) is a method of analyzing descriptive data using technology and knowledge in order to make informed business decisions [7]. The BI toolkit is used to collect, govern, and transform data. It allows internal and external stakeholders to communicate data, making decision-making easier. Business intelligence is the process of extracting useful information from data. Some of the things that BI can help you with are:

- Developing a better grasp of the marketplace
- Identifying new revenue streams
- Enhancing business procedures
- Keeping one step ahead of the competition

Cloud computing has been the most important facilitator of BI in recent years. The cloud has enabled organizations to process more data, from more sources, and more efficiently than they could before cloud technologies were introduced.

Data science vs. business intelligence: Understanding the differences between data science and business intelligence is beneficial. Understanding how they work together is also beneficial. It is not a question of choose one over the other [8]. It all boils down to picking the proper solution to

receive the information you need. Most of the time, this means combining data science and business intelligence. The simplest approach to distinguish between the two is to think about data science in terms of the future and business intelligence in terms of the past and present. Data science is concerned with predictive and prescriptive analysis, whereas business

Table 1.1 Comparison of data science and business intelligence.

Factor	Data science	Business intelligence
Concept	It is a discipline that employs mathematics, statistics, and other methods to uncover hidden patterns in data.	It is a collection of technology, applications, and processes that businesses employ to analyze business data.
Focus	It is concentrated on the future.	It concentrated the past and present.
Data	It can handle both structured and unstructured data.	It is mostly concerned with structured data.
Flexibility	Data science is more adaptable since data sources can be added as needed.	It is less flexible because data sources for business intelligence must be planned ahead of time.
Method	It employs the scientific process.	It employs the analytic method.
Complexity	In comparison to business intelligence, it is more sophisticated.	When compared to data science, it is a lot easier.
Expertise	Data scientist is its area of competence.	Its area of specialization is for business users.
Questions	It addresses the questions of what will happen and what might happen.	It is concerned with the question of what occurred.
Tools	SAS, BigML, MATLAB, Excel, and other tools are used.	Insight Squared Sales Analytics, Klipfolio, ThoughtSpot, Cyfe, TIBCO Spotfire, and more solutions are among them.