

READING MATERIALS – EDUCATIONAL PURPOSES ONLY

BIG DATA ANALYTICS

TABLE OF CONTENTS

Chapter 1 Introduction to Big Data Analytics	7
1.1 An introduction to Big Data Analytics	7
1.1.1 What is Analytics?	7
1.1.2 What is Big Data Analytics?	8
1.1.3 History and Evolution of Big Data Analytics	9
1.1.4 Values of Big Data Analytics	9
1.1.5 Uses of Big Data analytics across different industries	10
1.2.Big Data Analytics - Overview	11
1.3.Big Data Analytics - Data Life Cycle	11
1.3.1 Traditional Data Mining Life Cycle	11
1.3.2 Big Data Life Cycle	14
1.4.Big Data Analytics - Methodology	17
1.5.ANALYTICS PROCESS MODEL	17
1.6.ANALYTICAL MODEL REQUIREMENTS	19
1.7.1 Types of Data Sources	21
1.7.2 Sampling	22
1.7.3 Types of data elements	23
1.7.4 Visual data exploration and exploratory statistical analysis	24
1.7.5 Missing values	25
1.7.6 Outlier detection and treatment	26
1.7.7 Categorization	29
1.7.8 Variable selection	33
1.7.9 Segmentation	35
1.8. TYPES OF DATA	35
1.8.1 Datification: The new forms of data	38
1.8.2 The anatomy of Big Data	46
1.9. Big Data Tools and Techniques	49
1.9.1 Understanding Big Data Storage	49
1.9.2 A Generals Overview of High-Performance Architecture	49

1.9.3 HDFS	50
1.9.4 MapReduce and YARN	52
1.9.5 Expanding the Big Data Application Ecosystem	53
1.9.6 ZOOKEEPER	54
1.9.7 HBASE	54
1.9.8 HIVE	54
1.9.9 PIG	55
1.9.10 MAHOUT	56
1.9.11 Considerations	57
Chapter 2 Big Data Platforms - Hadoop	58
2.1 A Brief History of Hadoop	58
2.2 Big Data & Hadoop – Restaurant Analogy	59
2.3 What is Hadoop?	64
2.3.1 Hadoop-as-a-Solution	64
2.3.2 Hadoop Features	65
2.3.3 Hadoop Core Components	66
2.3.4 Hadoop Ecosystem	68
2.3.5 Examples of Hadoop - 5 Real-World Use Cases	69
2.3.6 Last.fm Case Study	70
2.3.7 Advantages of using Hadoop	71
2.3.8 Disadvantages of using Hadoop	72
2.3.9 Is Hadoop an efficient use of resources?	73
2.3.10 The business case for Hadoop	74
2.3.11 Another way to look at value	74
2.3.12 What does Hadoop replace?	75
2.3.13 Problems that Hadoop solves	75
2.4 Hadoop Installation	76
2.4.1 Java Installation	76
2.4.2 SSH Installation	77
2.4.3 Hadoop Installation	78

2.5 Hadoop Modules	81
2.5.1 HDFS	81
2.5.2 YARN	87
2.5.3 MapReduce	87
Chapter 3 Big Data Storage and Analytics	93
3.1 Big Data Storage Concepts	93
3.1.1 Clusters	93
3.1.2 File Systems and Distributed File Systems	94
3.1.3 NoSQL	96
3.1.4 Sharding	96
3.1.5 Replication	98
3.1.6 Sharding and Replication	103
3.1.7 CAP Theorem	106
3.1.8 ACID	110
3.1.9 BASE	114
3.1.10 Case Study Example	117
3.2 Apache Mahout	118
3.2.1 Mahout's Story	119
3.2.2 What is Machine Learning?	120
3.2.3 Mahout's Machine Learning Themes	123
3.2.3.1 Recommender engines	123
3.2.3.2 Clustering	131
3.2.3.3 Classification	137
3.2.4 Tackling large scale with Mahout and Hadoop	142
Chapter 4 Machine Learning, Streams and Database on Spark	144
4.1 Spark: Real Time Cluster Computing Framework	144
4.1.1 Real Time Analytics	144
4.1.2 Why Spark when Hadoop is already there?	146
4.1.3 What is Apache Spark?	147

4.1.4 Features of Apache Spark	148
4.1.5 Getting Started With Spark	150
4.1.6 Using Spark with Hadoop	151
4.1.7 Spark Components	152
4.1.8 Earthquake Detection using Spark	155
4.2 Sentiment Analysis Using Apache Spark	160
4.2.1 What is Streaming?	160
4.2.2 Why Spark Streaming?	160
4.2.3 Spark Streaming Overview	161
4.2.4 Spark Streaming Features	161
4.2.5 Spark Streaming Workflow	161
4.2.6 Spark Streaming Fundamentals	162
4.2.7 Use Case – Twitter Sentiment Analysis	166
4.3 Spark MLlib – Machine Learning Library Of Apache Spark	171
4.3.1 What is Machine Learning?	172
4.3.2 Spark MLlib Overview	173
4.3.3 Spark MLlib Tools	173
4.3.4 MLlib Algorithms	174
4.3.5 Use Case – Movie Recommendation System	176
4.4 Spark SQL Tutorial – Understanding Spark SQL With Examples	181
4.4.1 Why Spark SQL Came Into Picture?	182
4.4.2 Limitations with Hive	182
4.4.3 Spark SQL Overview	182
4.4.4 Spark SQL Libraries	183
4.4.5 Features Of Spark SQL	184
4.4.6 Querying Using Spark SQL	188
4.4.7 Creating Datasets	192
4.4.8 Adding Schema To RDDs	194
4.4.9 RDDs As Relations	197
4.4.10 Caching Tables In-Memory	199
4.4.11 Loading Data Programmatically	199

Big Data Analytics

4.4.12 JSON Datasets 201

4.4.13 Hive Tables 202

References 206

CHAPTER 1

INTRODUCTION TO BIG DATA ANALYTICS

1.1 AN INTRODUCTION TO BIG DATA ANALYTICS [1]

Many of you would have probably heard about Big data Analytics. Have you ever wondered what it is all about and how it can help us? Big data analytics can be defined as a process of examining large and varied data sets. We use advanced analytics techniques against the large data to uncover the hidden patterns, unknown correlations, market trends, customer preferences and other useful information. This helps the organizations to make informed decisions.

To understand Big Data Analytics you have to first understand What Analytics is?

1.1.1 What is Analytics? [1, 3]

Analytics is an encompassing and multidimensional field. It uses mathematics, statistics, predictive modeling and machine-learning techniques to find meaningful patterns and knowledge in recorded data.

Analytics is a term that is often used interchangeably with data science, data mining, knowledge discovery, and others. The distinction between all those is not clear cut. All of these terms essentially refer to extracting useful business patterns or mathematical decision models from a preprocessed data set. Different underlying techniques can be used for this purpose, stemming from a variety of different disciplines, such as:

- Statistics (e.g., linear and logistic regression)
- Machine learning (e.g., decision trees)
- Biology (e.g., neural networks, genetic algorithms, swarm intelligence)
- Kernel methods (e.g., support vector machines)

Basically, a distinction can be made between predictive and descriptive analytics. In predictive analytics, a target variable is typically available, which can either be categorical (e.g., churn or not, fraud or not) or continuous (e.g., customer lifetime value, loss given default). In descriptive analytics, no such target variable is available. Common examples here are association rules, sequence rules, and clustering. Figure 1.1 provides an example of a decision tree in a classification predictive analytics setting for predicting churn.

More than ever before, analytical models steer the strategic risk decisions of companies. For example, in a bank setting, the minimum equity and provisions a financial institution

holds are directly determined by, among other things, credit risk analytics, market risk analytics, operational risk analytics, fraud analytics, and insurance risk analytics. In this setting, analytical model errors directly affect profitability, solvency, shareholder value, the macroeconomic, and society as a whole. Hence, it is of the utmost importance that analytical models are developed in the most optimal way, taking into account various requirements that will be discussed in what follows.

Customer	Age	Recency	Frequency	Monetary	Churn
John	35	5	6	100	Yes
Sophie	18	10	2	150	No
Victor	38	28	8	20	No
Laura	44	12	4	280	Yes

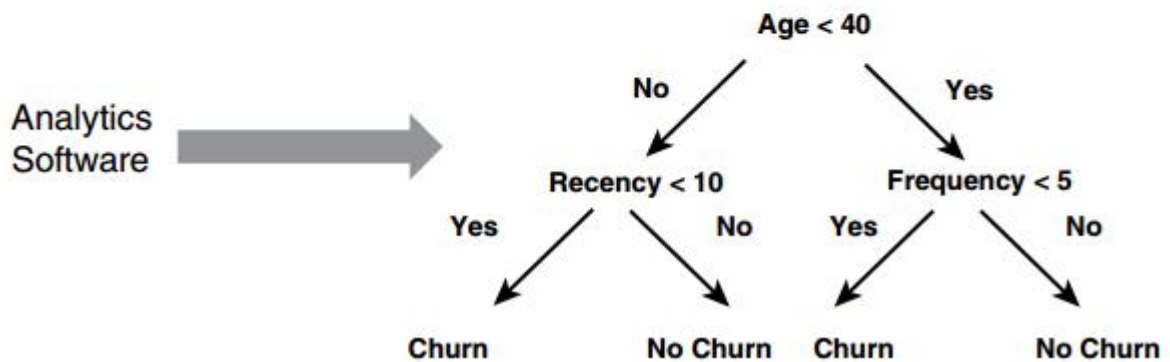


Figure 1.1 Example of Classification Predictive Analytics [3]

1.1.2 What is Big Data Analytics? [1]

As said before Big data analytics examines large amounts of data to uncover hidden patterns, correlations and other insights. With today's technology, it's possible to analyze your data and get answers from it immediately. Big Data Analytics helps you to understand your organization better. With the use of Big data analytics, one can make the informed decisions without blindly relying on guesses.

And it can help answer the following types of questions:

What actually happened?

How or why did it happen?

What's happening now?

What is likely to happen next?

1.1.3 History and Evolution of Big Data Analytics [1]

The concept of big data has been around for years; most organizations now understand that if they capture all the data that streams into their businesses, they can apply analytics and get significant value from it. But even in the 1950s, decades before anyone uttered the term “big data,” businesses were using basic analytics essentially numbers in a spreadsheet that were manually examined to uncover insights and trends.

The new benefits that big data analytics brings to the table, however, are speed and efficiency. Whereas a few years ago a business would have gathered information, run analytics and unearthed information that could be used for future decisions, today that business can identify insights for immediate decisions. The ability to work faster – and stay agile – gives organizations a competitive edge they didn't have before.

1.1.4 Values of Big Data Analytics [1]



Figure 1.2 Big Data Analytics

Big data analytics helps organizations harness their data and use it to identify new opportunities. That, in turn, leads to smarter business moves, more efficient operations, higher profits and happier customers. Here are the most important values of Big Data,

1. *Cost reduction:* Big data technologies such as Hadoop and cloud-based analytics bring significant cost advantages when it comes to storing large amounts of data – plus they can identify more efficient ways of doing business.

2. *Faster, better decision making:* With the speed of Hadoop and in-memory analytics, combined with the ability to analyze new sources of data, businesses are able to analyze information immediately – and make decisions based on what they’ve learned.
3. *New products and services:* With the ability to gauge customer needs and satisfaction through analytics comes the power to give customers what they want. Davenport points out that with big data analytics, more companies are creating new products to meet customers’ needs.

1.1.5 Uses of Big Data analytics across different industries[1]

Banking

Large amounts of information will be streaming in into banks, managing all this data and getting proper insights would be possible only with big data analytics. This is important to understand customers and boost their satisfaction, and also to minimize risk and fraud.

Government

When government agencies are able to harness and apply analytics to their big data, they gain significant ground when it comes to managing utilities, running agencies, dealing with traffic congestion or preventing crime.

Health Care

Patient records, Treatment plans, Prescription information. When it comes to health care, everything needs to be done quickly, accurately – and, in some cases, with enough transparency to satisfy stringent industry regulations. When big data is managed effectively, health care providers can uncover hidden insights that improve patient care.

Education

Educators armed with data-driven insight can make a significant impact on school systems, students, and curriculums. By analyzing big data, they can identify at-risk students, make sure students are making adequate progress, and can implement a better system for evaluation and support of teachers and principals.

Manufacturing

Armed with insight that big data can provide, manufacturers can boost quality and output while minimizing waste – processes that are key in today’s highly competitive market. More and more manufacturers are working in an analytics-based culture, which means they can solve problems faster and make more agile business decisions.