

# From Raw Data to Informed Decisions: Analyzing Amazon Book Reviews

Alberti A. • Ligari D. • Andreoli C.<sup>1</sup>

<sup>1</sup>Department of Computer Engineering - Data Science, University of Pavia, Italy  
Course of Data science and big data analytics

Github repository: <https://github.com/DavideLigari01/data-science-project>

---

## Abstract

**Keywords**— DNS reflection and amplification Attacks • Amplification factor • Ping • Wireshark • DIG • Mitigation measures

---

## Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Workplan</b>	<b>2</b>
<b>3</b>	<b>Discovery</b>	<b>2</b>
3.1	Team . . . . .	2
3.2	Tools . . . . .	2
3.3	Framing . . . . .	2
<b>4</b>	<b>Data Preparation</b>	<b>2</b>
4.1	Data Collection . . . . .	2
4.2	Hypothesis Generation . . . . .	2
4.3	Data Cleaning . . . . .	2
4.4	Data Aggregation . . . . .	2
4.5	MongoDB loading . . . . .	2
<b>5</b>	<b>Local Hypotheses Testing</b>	<b>2</b>
5.1	Hypothesis 1 . . . . .	2
5.2	Hypothesis 2 . . . . .	2
5.3	Hypothesis 3 . . . . .	2
5.4	Hypothesis 4 . . . . .	2
5.5	Hypothesis 5 . . . . .	2
5.6	Hypothesis 6 . . . . .	2
<b>6</b>	<b>Spark Hypotheses Testing</b>	<b>2</b>
<b>7</b>	<b>Helpfulness Prediction</b>	<b>2</b>
<b>8</b>	<b>Conclusions</b>	<b>2</b>

## 1 Introduction

In the age of digital commerce, customer reviews play a pivotal role in shaping product perception and influencing purchasing decisions. With the proliferation of online bookstores, Amazon has amassed an immense repository of book reviews spanning nearly two decades. These reviews contain valuable insights, sentiments, and trends that can unlock a treasure trove of information for authors, publishers, and book enthusiasts. This project embarks on a journey to harness the power of data, employing a comprehensive workflow to dissect and understand the vast collection of Amazon Books Reviews. Our mission is to develop a scalable solution that allow us to discover patterns, sentiment trends, and hidden correlations within the world of book reviews. We leverage cutting-edge tools and technologies, including Hadoop, Spark, MongoDB, and Python libraries such as Pandas and Scikit-learn. In this report, we embark on a detailed exploration of our project, delving into each stage of our workflow, from initial data discovery and preparation to feature extraction, model building, and rigorous evaluation.

## **2 Workplan**

## **3 Discovery**

### **3.1 Team**

### **3.2 Tools**

### **3.3 Framing**

## **4 Data Preparation**

### **4.1 Data Collection**

### **4.2 Hypothesis Generation**

### **4.3 Data Cleaning**

### **4.4 Data Aggregation**

### **4.5 MongoDB loading**

## **5 Local Hypotheses Testing**

### **5.1 Hypothesis 1**

### **5.2 Hypothesis 2**

### **5.3 Hypothesis 3**

### **5.4 Hypothesis 4**

### **5.5 Hypothesis 5**

### **5.6 Hypothesis 6**

## **6 Spark Hypotheses Testing**

## **7 Helpfulness Prediction**

## **8 Conclusions**