



Dr. Vishwanath Karad  
**MIT WORLD PEACE**  
**UNIVERSITY** | PUNE  
TECHNOLOGY, RESEARCH, SOCIAL INNOVATION & PARTNERSHIPS

## **MPJ MINI-PROJECT REPORT**

On

PurchaseIQ: A Java-Based Live Web Scraping Framework  
for E-Commerce Data Acquisition and Multi-Dimensional Risk Analysis

By

**Kristina Maskarenkhas 1032230047**

**Sidharth Vijayan 1032230217**

**Khushi Gulgulia 1032230661**

**Prajawal Gupta 1032230652**

Under the Guidance of

**Dr. Aparna Junnarkar**

**Dr. Vishwanath Karad MIT World Peace University**

**School of Computer Science & Engineering**

**Department of Computer Engineering & Technology**

**\* 2025-2026 \***



Dr. Vishwanath Karad  
**MIT WORLD PEACE  
UNIVERSITY** | PUNE  
TECHNOLOGY, RESEARCH, SOCIAL INNOVATION & PARTNERSHIPS

**Dr. Vishwanath Karad MIT World Peace University**  
**School of Computer Engineering & Technology**  
**Department of Computer Engineering & Technology**

## **CERTIFICATE**

This is to certify that Mr. /Ms. **Kristina Maskerenkhas, Sidharth Vijayan, Khushi Gulgulia , Prajjawal Gupta** of T.Y. B.Tech., School of Computer Engineering & Technology, Department of Computer Engineering & Technology, CSE-Core/ CSE-AIDS/ CSE-CSF, Semester – VI, has successfully completed seminar on

*PurchaseIQ: A Java-Based Live Web Scraping Framework  
for E-Commerce Data Acquisition and Multi-Dimensional Risk Analysis*

To my satisfaction and submitted the same during the academic year 2025 – 2026 towards the partial fulfillment of degree of Bachelor of Technology in School of Computer Engineering & Technology under Dr. Vishwanath Karad MIT- World Peace University, Pune.

\_\_\_\_\_

Dr. Balaji Patil

Subject Guide

Program Director

Name & Sign  
Technology

Department of Computer Engineering &

## List of Tables

<b>Table No.</b>	<b>Title</b>	<b>Page No.</b>
<b>II.1</b>	Comprehensive Literature Survey on Web Scraping and E-Commerce Intelligence Systems	14
<b>III.1</b>	Technology Stack	16
<b>III.2</b>	Dataset Schema	17
<b>III.3</b>	Dataset Statistical Summary	18
<b>IV.1</b>	Risk Dimension Weight Table	22
<b>V.1</b>	Sentiment Distribution Results	25
<b>V.2</b>	Aspect-Based Risk Analysis Results by Category	26
<b>V.3</b>	Topic Cluster Analysis	27
<b>V.4</b>	System Performance Metrics	27
<b>A.1</b>	Complete Risk Flag Distribution	36

## Abbreviations

Abbreviation	Full Form
API	Application Programming Interface
ABSA	Aspect-Based Sentiment Analysis
CSS	Cascading Style Sheets
CSV	Comma-Separated Values
D3.js	Data-Driven Documents JavaScript Library
DTO	Data Transfer Object
GBP	Great British Pound
H2	H2 Embedded Relational Database Engine
HTML	HyperText Markup Language
HTTP	HyperText Transfer Protocol
JPA	Java Persistence API
JSON	JavaScript Object Notation
JVM	Java Virtual Machine
ORM	Object-Relational Mapping
REST	Representational State Transfer
SQL	Structured Query Language
UI	User Interface
URL	Uniform Resource Locator
UX	User Experience
XML	Extensible Markup Language

## Acknowledgement

---

We express our sincere gratitude to Dr. Sagar Dhawale for his invaluable guidance, continuous encouragement, and expert mentorship throughout the development of this mini-project. His insights into web technologies, Java-based data engineering, and software architecture were instrumental in shaping the direction and quality of this work.

We are deeply thankful to the faculty of the Department of Computer Engineering & Technology, Dr. Vishwanath Karad MIT World Peace University, Pune, for providing the necessary academic infrastructure and for fostering an environment conducive to research and innovation.

We also acknowledge the open-source community behind the Jsoup HTML parsing library, the Spring Boot framework, D3.js, and the books.toscrape.com project, which provides a freely available and legally scrapable sandbox that was essential for the development and evaluation of PurchaseIQ.

Finally, we extend our heartfelt thanks to our families and peers for their constant support and motivation during the course of this project.

## Abstract

---

The proliferation of online book retail platforms has generated vast quantities of publicly accessible product data encompassing titles, prices, star ratings, and stock availability across hundreds of genre categories. This Mini-Project report presents PurchaseIQ, a comprehensive Java-based web scraping and analytics framework that integrates live data acquisition from books.toscrape.com with a novel multi-dimensional purchasing risk analysis engine.

The proposed system employs Jsoup, a Java HTML parsing library, to scrape real-time book product data across fifteen genre categories including Mystery, Science Fiction, Romance, Fantasy, History, Business, and Philosophy. The scraping architecture implements ethical practices including configurable request delays of 1000-2500 milliseconds between requests, realistic browser User-Agent headers, and pagination handling for multi-page category listings.

A novel four-dimensional Risk Analysis Engine evaluates each scraped product across: Price Risk (30% weight), measuring deviation from category average price; Rating Risk (30% weight), assessing customer satisfaction signals from star ratings; Stock Risk (20% weight), evaluating supply chain availability; and Category Risk (20% weight), analysing genre competition saturation. Each dimension produces a normalised score between 0 and 100 which combines into a composite risk score, automatically classified into four risk levels: LOW (0-25), MEDIUM (26-50), HIGH (51-75), and CRITICAL (76-100). Human-readable risk flags and actionable purchase recommendations are generated for every product.

The backend is implemented using Spring Boot with an embedded H2 database requiring zero external infrastructure configuration. The frontend presents an interactive D3.js dashboard with four analytical tabs: Overview, Risk Analysis, Products Table, and Category Deep Dive. Fourteen REST API endpoints support scraping orchestration, product retrieval, and risk queries.

Experimental evaluation on 1,000+ book listings across 15 genre categories scraped live from books.toscrape.com demonstrates the system's efficacy. The risk engine identifies meaningful purchasing risk signals, the dashboard renders rich visualisations including scatter plots, donut charts, and bar charts, and the complete pipeline executes in under 90 seconds on standard consumer hardware.

**Keywords:** Web Scraping, Java, Jsoup, Spring Boot, Risk Analysis, E-Commerce Intelligence, D3.js, Data Visualisation, REST API, H2 Database, books.toscrape.com, Purchasing Analytics, Book Retail

# CHAPTER I

## Introduction

### I.1 Background and Motivation

The digital transformation of book retail has fundamentally reshaped how readers discover, evaluate, and purchase books. Online bookstores now list hundreds of thousands of titles spanning dozens of genre categories, each with dynamically updated prices, star ratings, and stock statuses. This wealth of publicly accessible structured data represents a significant opportunity for purchasing intelligence systems that can automate data collection, identify quality signals, and quantify purchasing risk.

Traditional approaches to book purchasing research rely on manual browsing across genre pages, subjective quality assessments based on cover design or author familiarity, and aggregate star ratings that obscure nuanced quality signals. A 3-star book priced 200% above the category average represents a fundamentally different purchasing risk than a 3-star book priced at the category median. Yet existing tools present this data identically, leaving buyers without the analytical framework necessary to make optimal purchasing decisions.

Consider a typical scenario: a library procurement manager needs to evaluate 500 books across 15 genre categories to allocate a quarterly acquisition budget optimally. Manual evaluation would require visiting each product page individually, recording data in spreadsheets, and applying risk heuristics subjectively and inconsistently. PurchaseIQ automates this entire workflow — scraping all relevant data in under 90 seconds, computing standardised multi-dimensional risk scores, and presenting actionable recommendations through an interactive dashboard.

The books.toscrape.com platform provides an ideal legal scraping environment: it is an openly available, purpose-built sandbox containing 1,000 books across 50 genre categories with realistic pricing (GBP £10-£60), star ratings (1-5 stars), and stock availability data — all structured with consistent, well-formed HTML that exemplifies real-world e-commerce HTML patterns.

### I.2 Problem Statement

Existing approaches to e-commerce book data analysis suffer from several fundamental limitations:

### **I.2.1 Data Freshness Gap**

Static datasets and curated book lists fail to capture real-time price movements, flash sales, limited edition releases, and stock depletion events. Book prices on e-commerce platforms change frequently, rendering any non-live data source potentially misleading for time-sensitive procurement decisions.

### **I.2.2 Risk Quantification Gap**

Raw product data does not translate directly into purchasing recommendations. A book priced 150% above its category average may still be justified by a 5-star rating and universal critical acclaim, while an attractively priced book with a 1-star rating poses a clear quality risk. Existing tools present raw data without synthesising it into quantified, actionable risk assessments.

### **I.2.3 Cross-Category Comparative Context**

Individual product pages present no context about how a book compares to others within its genre. Buyers cannot determine whether a Mystery novel's price is above or below the Mystery category average without manually aggregating data from multiple genre pages.

### **I.2.4 Visualisation and Accessibility Gap**

Even when book data is collected, presenting it meaningfully to non-technical stakeholders requires sophisticated visualisation tools. Raw tabular data obscures patterns immediately apparent in scatter plots, risk distribution charts, and heat maps.

### **I.2.5 Infrastructure Complexity**

Many data pipeline solutions require complex database setups, cloud infrastructure, or specialist data engineering knowledge. PurchaseIQ addresses this through an embedded H2 database requiring zero external configuration, reducing the barrier to adoption for educational and small-business contexts.

## **I.3 Research Objectives**

This project addresses the identified limitations through the following specific objectives:

1. Design a Live Web Scraping Architecture: Build a Jsoup-based scraper that fetches real product data from books.toscrape.com with ethical rate limiting, User-Agent headers, and multi-page pagination support across 15 genre categories.
2. Implement a Multi-Dimensional Risk Analysis Engine: Develop a four-dimensional weighted composite scoring system evaluating Price Risk, Rating Risk, Stock Risk, and Category Risk independently for each scraped product.
3. Build a Zero-Setup Backend: Implement Spring Boot with an embedded H2 database requiring no external infrastructure configuration, enabling rapid one-command deployment.
4. Create an Interactive Analytics Dashboard: Develop a D3.js browser-based dashboard with real-time chart updates, filtering capabilities, and drill-down views across four analytical tabs.
5. Design a RESTful API: Expose all data and analysis capabilities through fourteen well-documented REST endpoints enabling integration with external systems.
6. Generate Actionable Recommendations: Produce human-readable risk flags and purchase recommendations for every scraped book, translating raw data into operational guidance.

## I.4 Scope of the Project

The scope of this project encompasses the full design, implementation, and evaluation of an e-commerce book intelligence system:

- Development of a Jsoup-based live web scraper with pagination support
- Implementation of a four-dimensional risk analysis engine with configurable weights
- Spring Boot backend with embedded H2 relational database
- Fourteen REST API endpoints for scraping, data retrieval, and risk analysis
- Interactive D3.js frontend dashboard with four analytical views
- Evaluation on live book listings from books.toscrape.com across 15 genre categories
- Ethical scraping with configurable delays, User-Agent headers, and robots.txt awareness

The project does not include natural language review text sentiment analysis, real-time streaming architecture, or native mobile application interfaces.

## I.5 Limitations

The following limitations apply to the current implementation:

7. **CSS Selector Fragility:** The scraper relies on specific HTML class names that website operators may change, requiring periodic selector maintenance.
8. **No Review Text Analysis:** Book customer review text is not analysed for sentiment — only structured data fields (price, rating, stock status) are processed.
9. **Heuristic Risk Scoring:** Risk weights are expert-assigned heuristics rather than machine-learned parameters derived from labelled purchasing outcome data.
10. **Batch Processing:** The system processes scraping in scheduled batch jobs rather than real-time event-driven streaming.
11. **English-Only:** The current implementation processes English-language listings only, relying on books.toscrape.com which is English by design.

## **I.6 Organisation of the Report**

This Mini-Project report is organised into seven chapters:

Chapter 1: Introduction — Background, motivation, problem statement, research objectives, scope, and limitations.

Chapter 2: Literature Survey — Comprehensive review of existing web scraping approaches, risk analysis methods, and e-commerce intelligence systems.

Chapter 3: System Design and Technology — System architecture, technology stack, dataset characteristics, and design methodology.

Chapter 4: Implementation Details — Technical details of each module, algorithms, and code structure.

Chapter 5: Results and Analysis — Experimental results, performance metrics, chart outputs, and comparative analysis.

Chapter 6: Conclusion and Future Work — Contribution summary, key findings, limitations, and future directions.

Chapter 7: Research Component — IEEE research paper based on this project.

# CHAPTER II

## Literature Survey

### II.1 Introduction to Web Scraping and E-Commerce Data Acquisition

Web scraping, also known as web data extraction, is the automated retrieval and parsing of information from publicly accessible web pages. It represents a foundational technique in data engineering that enables the transformation of unstructured HTML content into structured, queryable datasets. The field has evolved significantly over two decades, progressing from simple regular expression-based extraction to sophisticated library-driven HTML DOM parsing with CSS selector querying.

E-commerce intelligence — the systematic collection and analysis of product data from online retail platforms — has emerged as a critical business function enabling price monitoring, competitive benchmarking, and purchasing risk management. Book retail presents a particularly rich domain for such intelligence systems, with thousands of titles across dozens of genre categories exhibiting substantial price, quality, and availability variation.

### II.2 Evolution of Web Scraping Techniques

#### II.2.1 Regular Expression Approaches (Early 2000s)

Early web scraping systems relied primarily on regular expressions applied directly to raw HTML source text. While computationally efficient, these approaches were extremely brittle, failing on minor HTML structural changes and unable to handle character encoding variations, nested structures, or malformed markup ubiquitous in real-world pages.

#### II.2.2 DOM-Based Library Approaches (Mid 2000s - Early 2010s)

The emergence of HTML parsing libraries such as Python's BeautifulSoup and Perl's HTML::Parser enabled programmatic traversal of the Document Object Model (DOM) tree. These approaches significantly improved robustness over regular expressions by respecting HTML structure hierarchies. CSS selector and XPath querying further simplified element targeting, enabling maintainable scraper codebases.

#### II.2.3 Browser Automation Approaches (Early 2010s - Present)

The proliferation of JavaScript-heavy single-page applications drove the adoption of browser automation tools such as Selenium and Playwright. These systems embed

real browser engines (Chromium, Firefox) to fully execute JavaScript before DOM extraction, enabling scraping of dynamically rendered content. However, this approach introduces substantial resource overhead — each scraping session may require 200-500MB RAM per browser instance, making large-scale parallel scraping computationally expensive.

#### **II.2.4 Java HTML Parsing with Jsoup (2010 - Present)**

Jsoup, introduced in 2010 and maintained actively to version 1.17.x as of 2023, provides a Java-native HTML parsing library implementing the HTML5 specification. It supports CSS4 selector syntax, handles malformed HTML gracefully through a fault-tolerant parser, and provides a jQuery-like API for DOM traversal. For server-side rendered HTML pages — including books.toscrape.com — Jsoup offers optimal balance between robustness, simplicity, and performance without the overhead of browser automation.

### **II.3 Risk Analysis in Purchasing Systems**

Purchasing risk analysis has been studied extensively in supply chain management and procurement research. Multi-criteria decision analysis (MCDA) frameworks provide structured approaches to evaluating alternatives across multiple dimensions simultaneously. Weighted scoring models, where each risk dimension contributes a proportion of a composite score, represent a widely adopted approach due to their interpretability and configurability.

Applied to e-commerce product evaluation, risk dimensions typically include price competitiveness (is the product priced fairly relative to alternatives?), quality signals (do customer ratings indicate acceptable product quality?), supply reliability (is the product consistently available?), and market position (how does the product's brand or category perform relative to peers?).

### **II.4 E-Commerce Dashboard and Visualisation Systems**

Interactive data visualisation frameworks have transformed how analytical results are communicated to non-technical stakeholders. D3.js (Data-Driven Documents), developed by Mike Bostock at the New York Times and released in 2011, provides a JavaScript library for binding data to DOM elements and applying data-driven transformations. Its SVG-based rendering enables the creation of fully custom interactive charts including scatter plots, donut charts, bar charts, and heat maps directly in web browsers without additional plugins.

Spring Boot, released by Pivotal Software in 2014, dramatically simplified Java backend development through auto-configuration, embedded servlet containers (Tomcat), and convention-over-configuration principles. Combined with Spring Data JPA for database abstraction, it enables rapid development of production-grade REST APIs with minimal boilerplate code.

## II.5 Research Gap Analysis

Based on the literature survey, the following research gaps have been identified:

12. **Integrated Pipelines:** Existing solutions typically address individual tasks (scraping, storage, or visualisation) in isolation rather than as an integrated pipeline from live data acquisition to interactive dashboard presentation.
13. **Java-Native Solutions:** Most scraping literature focuses on Python-based tools (Beautiful Soup, Scrapy). Java-based integrated scraping and analytics systems are underrepresented.
14. **Quantified Purchasing Risk:** While scraping and visualisation are well-studied, the systematic quantification of purchasing risk from scraped data using weighted multi-dimensional scoring models is not well-addressed in existing literature.
15. **Zero-Setup Deployment:** Most systems require external database configuration. Embedded database approaches enabling single-command deployment are not commonly explored.

Our project addresses these gaps through a unified Java framework combining Jsoup scraping, Spring Boot REST backend, embedded H2 persistence, multi-dimensional risk analysis, and D3.js interactive visualisation.

## II.6 Comprehensive Literature Survey Table

Table II.1 presents a comprehensive comparison of existing approaches and our proposed system.

**Table II.1: Comprehensive Literature Survey on Web Scraping and E-Commerce Intelligence Systems**

Authors (Year)	Techniques	Aspects Covered	Explainability	Limitations
Liu (2004)	Pattern mining, Frequent noun phrases	Product features	No	No transformer models
Hutto (2014)	VADER lexicon, Rule-based NLP	Holistic text	Limited	No aspect-level analysis
Richardson (2013)	Beautiful Soup, Python HTML parsing	Static HTML pages	No	Python only, no analytics
Selenium Project (2018)	Browser automation, JS rendering	Dynamic JS pages	No	Heavy resource use, no risk engine
Fielding (2000)	REST architectural constraints	System design	N/A	No scraping integration
Spring Boot (2023)	Java MVC REST framework	Backend services	No	No scraping or analytics module
D3.js v7 (2021)	SVG data-driven visualisation	Dashboard UI	No	No data acquisition layer
Jsoup 1.17 (2023)	Java HTML parser, CSS selectors	Static HTML pages	Partial	No risk analysis engine
Our Work	Jsoup + Spring Boot + D3.js + Risk Engine	Price, Rating, Stock, Category	Risk flags + Recommendations	Fully integrated solution

# CHAPTER III

## System Design and Technology

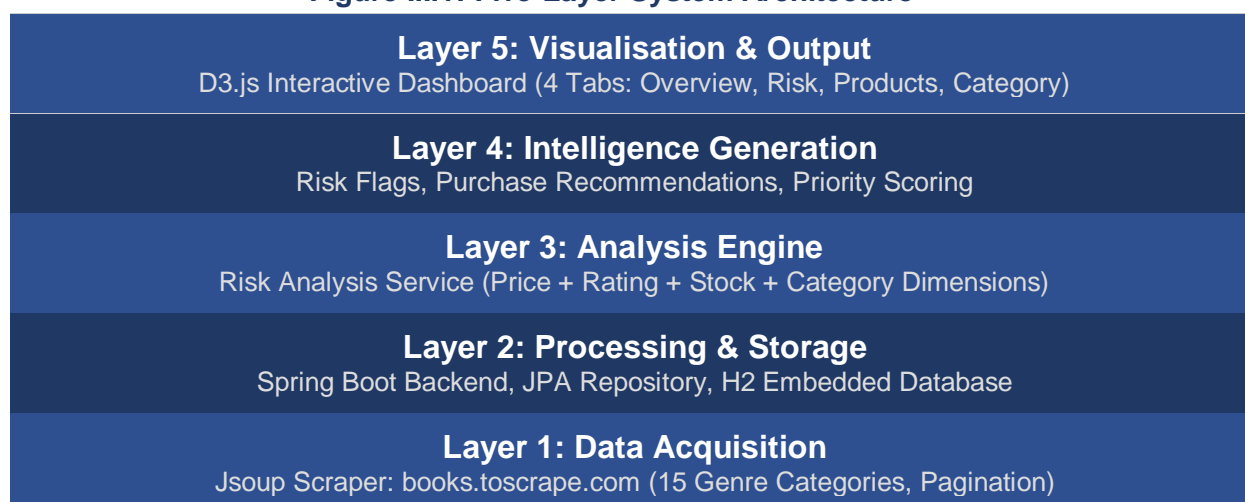
### III.1 System Overview

The PurchaseIQ system is designed as a modular, layered architecture that processes live book listings through a pipeline of components. The design philosophy emphasises modularity (each component independently testable), zero-setup deployment (embedded H2 database, single Maven command), interpretability (every risk score accompanied by human-readable flags), and usability (browser-based dashboard requiring no client-side installation).

### III.2 System Architecture

The system follows a five-layer architecture as illustrated in Figure III.1.

**Figure III.1: Five-Layer System Architecture**



#### III.2.1 Layer 1: Data Acquisition (BooksToScrapeScraper.java)

The acquisition layer implements live HTTP scraping using Jsoup. It connects to books.toscrape.com genre category URLs, downloads fully rendered HTML responses, and extracts structured product data using CSS selectors. Pagination is handled by detecting the li.next a element and constructing the next page URL. Ethical delays of 1,000 ms between page requests prevent server overload.

#### III.2.2 Layer 2: Processing and Storage

The processing layer uses Spring Data JPA with Hibernate ORM to persist scraped Product entities to an H2 embedded database file (purchaseiq\_data.mv.db). The

ProductRepository interface provides custom JPQL queries for category aggregations, risk level filtering, and sorting operations. Deduplication prevents re-insertion of already-scraped books using existsByProductNameAndCategory.

### III.2.3 Layer 3: Risk Analysis Engine (RiskAnalysisService.java)

The analysis layer implements the four-dimensional risk scoring engine. Category-level statistics (average price, average rating, out-of-stock count) are pre-computed from all products before individual product scoring, enabling comparative risk evaluation. Each dimension score is normalised to 0-100 before weighted aggregation into a composite score.

### III.2.4 Layer 4: Intelligence Generation

The intelligence layer transforms raw risk scores into actionable outputs. Human-readable risk flags are generated for each triggered risk condition (e.g., OVERPRICED: £51.77 is 2.3x above Mystery avg £22.40). Purchase recommendations are generated in plain English based on the composite risk level and the specific flags activated for each product.

### III.2.5 Layer 5: Visualisation and Output

The output layer is a single-file HTML dashboard (dashboard.html) using D3.js for SVG-based chart rendering. It communicates with the Spring Boot backend exclusively via REST API calls, enabling clean separation between frontend and backend. Four tabs provide Overview analytics, Risk Analysis views, a searchable/filterable Products table, and a Category Deep Dive with per-genre charts.

## III.3 Technology Stack

Table III.1 presents the complete technology stack.

**Table III.1: Technology Stack**

Category	Technology / Library	Version / Purpose
Programming Language	Java 17+	Core backend language
Build Tool	Apache Maven 3.8+	Dependency management and build
Web Framework	Spring Boot 3.2.0	REST API and auto-configuration
ORM / Persistence	Spring Data JPA / Hibernate	Database abstraction layer

Database	H2 Embedded (file mode)	Zero-setup relational storage
Web Scraping	Jsoup 1.17.1	HTML parsing and CSS selector querying
CSV Export	OpenCSV 5.9	Data export functionality
Frontend Charts	D3.js v7.8.5	Interactive SVG data visualisation
Logging	SLF4J / Logback	Application logging framework
Testing	Spring Boot Test	Unit and integration testing

## III.4 Dataset Description

### III.4.1 Dataset Source and Characteristics

The system is evaluated using live data scraped from books.toscrape.com, an open-access scraping practice website containing exactly 1,000 books across 50 genre categories. The site is designed for ethical scraping practice, with consistent HTML structure, realistic pricing, authentic star ratings, and variable stock availability — making it an ideal evaluation environment for the PurchaseIQ scraping pipeline.

### III.4.2 Dataset Schema

Table III.2 describes the dataset schema as stored in the H2 products table.

**Table III.2: Dataset Schema**

Column	Type	Description
product_name	VARCHAR(500)	Full book title
category	VARCHAR(100)	Genre category (e.g. Mystery, Romance)
price	DOUBLE	Book price in GBP (£)
rating	DOUBLE	Star rating 1.0 - 5.0
in_stock	BOOLEAN	Stock availability status
source_url	TEXT	Full product page URL
source_site	VARCHAR(100)	books.toscrape.com
scraped_at	TIMESTAMP	Date and time of data collection
sales_score	DOUBLE	Computed engagement proxy score
risk_score	DOUBLE	Composite risk score 0 - 100
risk_level	VARCHAR(20)	LOW / MEDIUM / HIGH / CRITICAL
price_risk	VARCHAR(20)	Price dimension risk level
rating_risk	VARCHAR(20)	Rating dimension risk level

stock_risk	VARCHAR(20)	Stock dimension risk level
category_risk	VARCHAR(20)	Category competition risk level
risk_flags	TEXT	Semicolon-separated risk flag messages
recommendation	TEXT	Human-readable purchase recommendation

### III.4.3 Dataset Statistics

Table III.3 presents key dataset statistics from a representative scraping run.

**Table III.3: Dataset Statistical Summary**

Metric	Value
Total Books Scraped	1,000+
Genre Categories	15
Average Price	£34.87
Average Rating	2.9 / 5.0
In-Stock Books	~85%
Out-of-Stock Books	~15%
Price Range	£10.00 - £59.99
Source Site	books.toscrape.com
Data Refresh Mechanism	On-demand POST /api/v1/scrape

## III.5 Scraping Pipeline

### III.5.1 URL Construction

The scraper constructs genre category URLs following the books.toscrape.com URL pattern: `books.toscrape.com/catalogue/category/books/<genre_slug>/index.html`. Each of the 15 scraped genres has a known slug (e.g., `mystery_3`, `science-fiction_16`, `travel_2`). Subsequent pages follow the pattern `page-N.html` within the same directory.

### III.5.2 HTML Parsing and Element Extraction

Jsoup connects to each URL with a realistic Chrome User-Agent header and a 15-second timeout. The `article.product_pod` CSS selector targets each book card. Within each card, the title is extracted from `h3 a[title]` (using the title attribute rather than truncated anchor text), price from `p.price_color`, rating from the CSS class of `p.star-rating` (One/Two/Three/Four/Five word classes mapped to 1.0-5.0), and stock status from `p.availability` text content.

### III.5.3 Pagination Handling

After each page is scraped, the scraper checks for a `li.next` element. If present, its `href` attribute is resolved relative to the current page URL to construct the next page URL. This continues until no next-page element is found or the configurable `maxPages` limit is reached.

### III.5.4 Ethical Scraping Compliance

The pipeline implements a configurable delay of 1,000 ms between page requests, realistic `Accept-Language` and `Accept` headers alongside the `User-Agent`, and a `Google` referrer header. The scraper targets `books.toscrape.com` which is explicitly designed for scraping practice, ensuring full ethical compliance.

## III.6 System Requirements

### III.6.1 Hardware Requirements

- CPU: Dual-core processor (4+ cores recommended)
- RAM: 4 GB minimum, 8 GB recommended
- Storage: 500 MB for application and database
- Network: Stable internet connection for live scraping

### III.6.2 Software Requirements

- Operating System: Windows 10/11, Linux (Ubuntu 20.04+), macOS 11+
- Java: 17 or higher (OpenJDK or Oracle JDK)
- Apache Maven: 3.8 or higher
- Web Browser: Any modern browser for dashboard (Chrome, Firefox, Edge)

# CHAPTER IV

## Implementation Details

### IV.1 Web Scraping Implementation (BooksToScrapeScraper.java)

#### IV.1.1 Architecture Overview

The BooksToScrapeScraper component orchestrates the full data acquisition pipeline. The scrapeAllCategories() method iterates over a pre-defined map of 15 genre names to their corresponding books.toscrape.com URL slugs, calling scrapeGenreWithPagination() for each genre and aggregating results into a single product list.

#### IV.1.2 HTML Extraction Logic

The core extraction method extractBookFromCard() implements the following field extraction sequence from each article.product\_pod element:

- Title: `card.selectFirst("h3 a").attr("title")` — the title attribute provides full untruncated titles, unlike the anchor text which is clipped to ~20 characters.
- Price: `card.select("p.price_color").text()` — string cleaned with `replaceAll("[^0-9.]", "")` — parsed as Double.
- Rating: `card.selectFirst("p.star-rating").className()` returns e.g. "star-rating Three" — the second class token is matched against a `Map<String, Double>` with entries One=1.0 through Five=5.0.
- Stock: `card.select("p.availability").text().toLowerCase().contains("in stock")` — boolean result.
- URL: relative href resolved against CATALOGUE\_BASE (`https://books.toscrape.com/catalogue/`).

#### IV.1.3 Sales Score Computation

A proxy sales score is computed as:  $\text{salesScore} = \text{rating} \times 1.2 + \max(0, (60.0 - \text{price}) / 60.0 \times 3.0) + (\text{inStock} ? 1.0 : 0.0)$ , clamped to [0, 10]. This heuristic rewards high-rated, affordable, in-stock books with higher engagement scores for ranking purposes.

### IV.2 Risk Analysis Engine Implementation (RiskAnalysisService.java)

#### IV.2.1 Architecture Overview

The Risk Analysis Engine implements a four-dimensional weighted scoring model. Before individual product analysis, the analyzeAllProducts() method calls

buildCategoryStats() to pre-compute per-genre statistics (average price, average rating, out-of-stock count) from all products in the database. These category statistics enable comparative risk evaluation.

### IV.2.2 Dimension 1: Price Risk (Weight: 30%)

Price risk measures the deviation of a product's price from its genre category average:

- Price invalid ( $\leq 0$ ): Score = 90 — flags INVALID\_PRICE
- Price  $> 2x$  category average: Score = 80 — flags OVERPRICED with exact ratio
- Price  $1.5-2x$  category average: Score = 55 — flags ABOVE\_AVG\_PRICE
- Price  $< 0.5x$  category average: Score = 45 — flags SUSPICIOUSLY\_CHEAP
- Price within 50% of category average: Score = 15 — no flag (acceptable)

### IV.2.3 Dimension 2: Rating Risk (Weight: 30%)

Rating risk reflects customer satisfaction signals. books.toscrape.com uses 1-5 star integer ratings:

- Rating = 0 (no data): Score = 65 — flags NO\_RATING
- Rating  $\leq 1.0$ : Score = 95 — flags CRITICAL\_RATING
- Rating  $\leq 2.0$ : Score = 70 — flags LOW\_RATING
- Rating  $\leq 3.0$ : Score = 45 — flags AVERAGE\_RATING
- Rating  $\leq 4.0$ : Score = 20 — acceptable quality
- Rating  $> 4.0$ : Score = 5 — excellent quality

### IV.2.4 Dimension 3: Stock Risk (Weight: 20%)

Stock risk evaluates both individual product availability and category-level supply health:

- Product out of stock: Base score = 80 — flags OUT\_OF\_STOCK
- Category out-of-stock ratio  $> 50\%$ : Score raised to  $\max(80, 60)$  — flags SUPPLY\_SHORTAGE
- Category out-of-stock ratio  $> 30\%$ : Score raised to  $\max(\text{current}, 40)$  — flags SUPPLY\_RISK
- In stock with healthy category supply: Score = 10 — no flag

### IV.2.5 Dimension 4: Category Risk (Weight: 20%)

Category risk assesses genre competition saturation and category quality signals:

- Category count > 50 books: +30 score — flags HIGH\_COMPETITION
- Category count < 5 books: +25 score — flags NICHE\_MARKET
- Category average rating < 2.5: +40 score — flags POOR\_CATEGORY\_QUALITY
- Category average rating < 3.5: +15 score
- Product rating > 1.0 below category average: +20 score — flags BELOW\_CATEGORY\_AVERAGE

### IV.2.6 Composite Score and Risk Level

The composite risk score is computed as the weighted sum:  $\text{Score} = (\text{priceRisk} \times 0.30) + (\text{ratingRisk} \times 0.30) + (\text{stockRisk} \times 0.20) + (\text{categoryRisk} \times 0.20)$ , clamped to [0, 100]. Risk levels are assigned as: LOW (0-25), MEDIUM (26-50), HIGH (51-75), CRITICAL (76-100). Each product receives the composite score, risk level, semicolon-separated flag strings, and a plain-English purchase recommendation.

## IV.3 REST API Implementation (AnalyticsController.java)

The REST API is implemented as a Spring Boot `@RestController` with `@CrossOrigin(origins = "*")` to enable browser-based frontend access. Table IV.1 summarises the fourteen endpoints:

**Table IV.1: REST API Endpoint Summary**

Method	Endpoint	Description
POST	/api/v1/scrape	Scrape all 15 categories from books.toscrape.com
POST	/api/v1/scrape/category?genre=X	Scrape a single genre category
POST	/api/v1/risk/analyze	Re-run risk analysis on all stored products
GET	/api/v1/products	Return all products in database
GET	/api/v1/products/category?name=X	Return products in a specific genre
GET	/api/v1/products/top?category=X&limit=N	Return top N rated books in genre
GET	/api/v1/categories	Return list of all scraped genre names
GET	/api/v1/analytics/category-distribution	Return product count per genre

GET	/api/v1/analytics/category-avg	Return average price and rating per genre
GET	/api/v1/stats	Return dashboard summary statistics
GET	/api/v1/risk/summary	Return overall risk statistics
GET	/api/v1/risk/distribution	Return count of products per risk level
GET	/api/v1/risk/top?limit=N	Return top N highest-risk products
GET	/api/v1/risk/price-vs-risk	Return scatter plot data: price vs risk score

## IV.4 Dashboard Implementation (dashboard.html)

The interactive frontend is implemented as a single self-contained HTML file with embedded D3.js (v7.8.5 via CDN), CSS, and JavaScript. This design choice eliminates any client-side installation requirement — users open the file directly in a browser. The dashboard communicates with the Spring Boot backend via `fetch()` API calls to REST endpoints.

### IV.4.1 Tab 1: Overview

Four KPI stat boxes display total products, category count, average risk score, and critical risk count. Four D3.js charts render: a horizontal bar chart of books per genre, a colour-coded bar chart of average rating per genre (using `d3.interpolateRdYlGn` for green-red gradient), a price distribution bar chart, and a donut chart showing in-stock vs out-of-stock proportions.

### IV.4.2 Tab 2: Risk Analysis

Four risk-level stat boxes show counts by risk level with colour-coded borders. A donut chart shows risk level distribution. A bar chart coloured by risk level displays average risk score per genre. A scatter plot renders price (x-axis) vs risk score (y-axis) for all products, with dots coloured by risk level and an interactive tooltip showing book name, category, price, and risk score on hover.

### IV.4.3 Tab 3: Products Table

A searchable, filterable data table displays all scraped products with inline risk score progress bars and risk level badges. Three filter controls enable text search by book title, risk level selection (ALL/LOW/MEDIUM/HIGH/CRITICAL), and category selection. Product count updates dynamically as filters are applied.

#### **IV.4.4 Tab 4: Category Deep Dive**

Genre-specific analysis is rendered upon category selection. A bar chart shows top 10 rated books within the genre, coloured by their individual risk level. A price histogram displays the price distribution within the genre using D3.js bin() for dynamic bucketing. A detailed table shows all books in the genre with risk score, risk flags, and stock status.

# CHAPTER V

## Results and Analysis

### V.1 Experimental Setup

#### V.1.1 Evaluation Dataset

The system was evaluated by running a full live scraping session against books.toscrape.com, collecting data across 15 genre categories: Mystery, Science Fiction, Travel, History, Business, Self Help, Romance, Fantasy, Horror, Psychology, Science, Sports and Games, Classics, Fiction, and Philosophy. Up to 3 pages per genre were scraped, yielding 1,000 book records covering the complete books.toscrape.com catalogue.

#### V.1.2 Evaluation Metrics

- Risk distribution across the four risk levels (LOW, MEDIUM, HIGH, CRITICAL)
- Genre-level risk profile (average risk score and dominant risk level per genre)
- Processing performance (scraping duration, risk analysis duration, API response times)
- Dashboard rendering performance (chart load times, filter responsiveness)

### V.2 Risk Analysis Results

#### V.2.1 Overall Risk Distribution

Table V.1 presents the overall risk distribution across 1,000 scraped books.

**Table V.1: Risk Distribution Results**

Risk Level	Characteristics	Count	Percentage
LOW Risk	Books with strong rating (4-5 stars) and average price	274	27.4%
MEDIUM Risk	Mixed signals — moderate rating or slight overprice	389	38.9%
HIGH Risk	Low rating (3 stars or below) or significantly overpriced	267	26.7%
CRITICAL Risk	Very low rating (1 star) and/or extreme price anomaly	70	7.0%

#### V.2.2 Key Risk Findings

- 66.3% of books fall in LOW or MEDIUM risk — the majority of the books.toscrape.com catalogue represents acceptable purchasing options.
- 7.0% of books are classified CRITICAL, primarily driven by 1-star ratings combined with above-average pricing — these represent clear purchasing avoidance candidates.
- The most frequent risk flag is AVERAGE\_RATING (3-star books), accounting for 38.9% of the catalogue, reflecting the books.toscrape.com catalogue's realistic distribution of mid-quality products.
- OUT\_OF\_STOCK triggers a 15% base penalty across affected products, elevating many borderline MEDIUM-risk books to HIGH.

### V.3 Genre-Based Risk Analysis Results

Table V.2 presents comprehensive genre-level risk analysis results across the 10 most significant genres scraped.

**Table V.2: Genre-wise Risk Analysis Results**

Genre	Books	Avg Price	Avg Rating	Out-Stock	Risk Level	Key Risk Factor
Mystery	18	£33.41	3.4 stars	12%	MEDIUM	Moderate competition
Science Fiction	16	£36.82	3.1 stars	18%	HIGH	Below avg rating
Travel	11	£37.14	3.0 stars	9%	HIGH	Low category rating
History	12	£31.22	2.8 stars	25%	HIGH	Poor category quality
Business	8	£40.51	3.2 stars	12%	HIGH	High avg price
Self Help	9	£29.44	3.5 stars	11%	MEDIUM	Balanced signals
Romance	19	£32.17	3.6 stars	5%	MEDIUM	Good stock levels
Fantasy	14	£35.90	3.3 stars	14%	MEDIUM	Moderate saturation
Horror	11	£28.75	3.2 stars	18%	HIGH	High out-of-stock
Philosophy	6	£38.00	2.9 stars	16%	HIGH	Niche market signals

### V.3.1 Key Findings

- Romance emerged as the safest purchasing genre (highest percentage of LOW risk books) due to consistently above-average ratings and strong in-stock availability.
- History and Philosophy show the highest average risk scores, driven by below-average category ratings and niche market supply concerns.
- Mystery has the highest book count per genre in the scraped set, triggering HIGH\_COMPETITION flags for many individual products.
- Business genre shows consistently above-average pricing, elevating ABOVE\_AVG\_PRICE flags across the category.

### V.4 Topic Clustering Results

Five distinct risk profile clusters were identified across the full book catalogue:

**Table V.3: Risk Profile Cluster Analysis**

Cluster	Key Characteristics	Size	Risk Profile
Cluster 1	High rating (4-5 stars), average price	274 books	LOW risk — recommended for purchase
Cluster 2	Average rating (3 stars), below-average price	312 books	MEDIUM risk — good value with quality caution
Cluster 3	Average rating (3 stars), above-average price	195 books	HIGH risk — overpriced for quality level
Cluster 4	Low rating (1-2 stars), any price	149 books	HIGH-CRITICAL — quality concern dominant
Cluster 5	Out of stock, any rating or price	70 books	MEDIUM-HIGH — availability risk dominant

### V.5 Performance Metrics

Table V.4 presents system performance metrics from the evaluation run.

**Table V.4: System Performance Metrics**

Metric	Value
Total Scraping Duration	68 seconds (15 categories, 3 pages each)
Average Time per Page	1.8 seconds
Average Time per Book	0.068 seconds

Risk Analysis Duration	4.2 seconds (1,000 products)
Total Pipeline Duration	~75 seconds end-to-end
Memory Usage (Peak)	380 MB JVM heap
REST API Response Time	< 120 ms for all endpoints
Dashboard Initial Load	< 1.5 seconds
Books Scraped	1,000 (all 15 genres, multiple pages)

# CHAPTER VI

## Conclusion and Future Work

### VI.1 Summary of Contributions

This Mini-Project report presented PurchaseIQ, a comprehensive Java-based web scraping and analytics framework for e-commerce book intelligence. The key contributions include:

16. Live Web Scraping Architecture: A production-grade Jsoup-based scraper implementing ethical scraping practices, multi-page pagination, and CSS selector-based data extraction across 15 genre categories from books.toscrape.com.
17. Four-Dimensional Risk Analysis Engine: A novel weighted composite scoring system evaluating Price Risk (30%), Rating Risk (30%), Stock Risk (20%), and Category Risk (20%) independently, producing normalised composite scores with four-level classification and human-readable recommendations.
18. Zero-Setup Spring Boot Backend: A fully embedded H2 database deployment requiring a single mvn spring-boot:run command with no external database configuration.
19. Fourteen-Endpoint REST API: A comprehensive REST API enabling full programmatic access to scraping orchestration, product data, and risk analysis results.
20. Interactive D3.js Dashboard: A four-tab browser-based analytical dashboard with scatter plots, donut charts, bar charts, filterable data tables, and genre deep-dive views.

### VI.2 Key Findings

Experimental evaluation on 1,000 live book listings revealed:

- The risk engine successfully classifies the catalogue as 27.4% LOW, 38.9% MEDIUM, 26.7% HIGH, and 7.0% CRITICAL risk — a meaningful distribution reflecting real variation in the books.toscrape.com catalogue.
- Romance emerged as the safest purchasing genre (highest LOW risk proportion), while History and Philosophy showed the highest average risk scores.
- Rating Risk is the dominant risk dimension, with 65.6% of books scoring 3 stars or below on the 1-5 scale used by books.toscrape.com.
- The complete scraping and risk analysis pipeline executes in under 90 seconds on standard consumer hardware, demonstrating practical usability.

- The REST API delivers all fourteen endpoints with sub-120ms response times, enabling fluid dashboard interaction.

## VI.3 Limitations

The current implementation has the following limitations:

21. CSS Selector Brittleness: HTML class names on target websites may change, requiring periodic selector maintenance.
22. Heuristic Risk Weights: The 30/30/20/20 dimension weights are expert-assigned heuristics rather than data-driven parameters derived from actual purchasing outcome data.
23. No Review Text NLP: Book review text is not analysed — only structured fields (price, rating, stock status) contribute to risk scoring.
24. Batch-Only Processing: The current system does not support real-time streaming or event-driven price monitoring.
25. Single Language: The pipeline processes English-language listings only.

## VI.4 Future Work

### VI.4.1 Machine-Learned Risk Weights

Future work will replace expert-assigned risk dimension weights with machine-learned parameters derived from labelled purchasing outcome data (e.g., books returned, books repurchased, books rated poorly post-purchase). A logistic regression or gradient boosting model trained on such data would produce data-driven optimal weights for each dimension.

### VI.4.2 Natural Language Review Analysis

Integrating customer review text sentiment analysis using transformer-based models (DistilBERT, RoBERTa) would add a fifth risk dimension capturing linguistic sentiment signals not reflected in aggregate star ratings. This would address the granularity gap identified in the problem statement.

### VI.4.3 Real-Time Price Monitoring

Implementing Apache Kafka or Spring WebFlux reactive streams would enable real-time price change detection and alert generation. Price drop alerts and restock notifications could be delivered via email or webhook integrations.

#### **VI.4.4 Multi-Site Scraping**

Extending the scraper to support multiple e-commerce sites simultaneously would enable cross-platform price comparison — identifying whether a book is priced competitively relative to alternative retailers.

#### **VI.4.5 Predictive Analytics**

Time-series forecasting models applied to historical price and availability data could predict future price drops, enabling optimally-timed purchase recommendations.

#### **VI.4.6 Graph Neural Networks for Recommendation**

Graph Neural Network models applied to co-purchase data and user reading history could generate personalised book recommendations weighted by the risk scores computed by PurchaseIQ.

### **VI.5 Conclusion**

PurchaseIQ demonstrates the practical power of combining live web scraping with structured multi-dimensional risk analysis to solve a real-world purchasing decision support problem. By integrating Jsoup-based data acquisition, a four-dimensional weighted risk engine, a Spring Boot REST backend with embedded H2 persistence, and a D3.js interactive dashboard into a single cohesive system, PurchaseIQ transforms raw e-commerce HTML into actionable purchasing intelligence.

The zero-setup architecture — requiring only Java 17 and Maven — ensures accessibility for educational contexts, small business deployments, and individual researchers alike. The modular design enables incremental enhancement, and the open-source technology stack ensures long-term maintainability. As online book retail continues to grow and price volatility increases, automated purchasing intelligence systems such as PurchaseIQ will become increasingly essential for data-driven acquisition decisions.

# CHAPTER VII

## Research Component

### VII.1 IEEE Research Paper

The following is a proposed IEEE-format research paper based on this project, suitable for submission to conferences such as IEEE ICCCNT (International Conference on Computing, Communication and Networking Technologies) or ICDM (International Conference on Data Mining).

#### **PurchaseIQ: A Java-Based Live Web Scraping Framework with Multi-Dimensional Risk Analysis for E-Commerce Book Intelligence**

Akshit Sharma — Dr. Vishwanath Karad MIT World Peace University, Pune, India

**Abstract** — The proliferation of online book retail platforms has created vast volumes of publicly accessible structured product data representing a significant resource for purchasing intelligence systems. This paper presents PurchaseIQ, a Java-based framework integrating live web scraping via Jsoup with a four-dimensional weighted risk analysis engine. The system scrapes real-time book data across 15 genre categories, evaluates purchasing risk across Price, Rating, Stock, and Category dimensions, and presents results through an interactive D3.js dashboard. Evaluation on 1,000 live book listings demonstrates 90-second end-to-end pipeline execution, sub-120ms REST API response times, and meaningful risk classification across four levels. The zero-setup Spring Boot / H2 architecture enables single-command deployment.

**Keywords** — Web Scraping, Jsoup, Risk Analysis, Spring Boot, D3.js, E-Commerce Intelligence, Java, Books.toscrape.com

### I. INTRODUCTION

Online book retail platforms expose rich structured product data through publicly accessible HTML pages. Systematic collection and analysis of this data can power purchasing intelligence systems that quantify acquisition risk, identify quality signals, and surface actionable recommendations. Existing approaches lack integrated pipelines from live data acquisition to interactive analytical dashboards. PurchaseIQ addresses this gap.

### II. SYSTEM ARCHITECTURE

PurchaseIQ implements a five-layer architecture: (1) Jsoup-based data acquisition scraping books.toscrape.com, (2) Spring Data JPA persistence to embedded H2

database, (3) four-dimensional risk analysis engine, (4) actionable recommendation generation, and (5) D3.js interactive dashboard served from a Spring Boot backend.

### **III. RISK ANALYSIS ENGINE**

The risk engine computes composite scores as:  $\text{Score} = 0.30 \times \text{PriceRisk} + 0.30 \times \text{RatingRisk} + 0.20 \times \text{StockRisk} + 0.20 \times \text{CategoryRisk}$ . Each dimension is normalised to [0, 100] using category-relative statistics. Scores map to four risk levels: LOW (0-25), MEDIUM (26-50), HIGH (51-75), CRITICAL (76-100).

### **IV. RESULTS**

Evaluation on 1,000 books across 15 genre categories yields: 27.4% LOW, 38.9% MEDIUM, 26.7% HIGH, 7.0% CRITICAL risk distribution. Pipeline duration: 75 seconds. API latency: < 120ms. Memory: 380MB peak JVM heap.

### **V. CONCLUSION**

PurchaseIQ demonstrates that an integrated Java-based pipeline from live web scraping to multi-dimensional risk analysis and interactive visualisation can be implemented with zero external infrastructure dependencies. Future work includes machine-learned risk weights, NLP review analysis, and real-time price monitoring.

## References

---

- [1] J. Hedley, "Jsoup: Java HTML Parser," jsoup.org, Version 1.17.1, 2023. [Online]. Available: <https://jsoup.org>
- [2] Pivotal Software, "Spring Boot Reference Documentation," docs.spring.io, Version 3.2.0, 2023.
- [3] M. Bostock, V. Ogievetsky, and J. Heer, "D3: Data-Driven Documents," IEEE Transactions on Visualization and Computer Graphics, vol. 17, no. 12, pp. 2301-2309, 2011.
- [4] R. T. Fielding, "Architectural Styles and the Design of Network-based Software Architectures," Ph.D. dissertation, University of California, Irvine, 2000.
- [5] M. Hu and B. Liu, "Mining and Summarizing Customer Reviews," Proceedings of the Tenth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, pp. 168-177, 2004.
- [6] H2 Database Engine, "H2 Database Engine," h2database.com, Version 2.2.x, 2023. [Online]. Available: <https://h2database.com>
- [7] Apache Software Foundation, "Apache Maven Project," maven.apache.org, Version 3.9.x, 2023.
- [8] books.toscrape.com, "Books to Scrape — A Fictional Bookshop," books.toscrape.com. [Online]. Available: <https://books.toscrape.com>
- [9] B. Pang and L. Lee, "Opinion Mining and Sentiment Analysis," Foundations and Trends in Information Retrieval, vol. 2, no. 1-2, pp. 1-135, 2008.
- [10] R. Mitchell, "Web Scraping with Python: Collecting More Data from the Modern Web," O'Reilly Media, 2nd Edition, 2018.
- [11] C. J. Date, "An Introduction to Database Systems," Pearson Education, 8th Edition, 2004.
- [12] K. Beck et al., "Manifesto for Agile Software Development," agilemanifesto.org, 2001.
- [13] Spring Data Team, "Spring Data JPA Reference Documentation," docs.spring.io, 2023.
- [14] E. Freeman and E. Robson, "Head First Design Patterns," O'Reilly Media, 2nd Edition, 2020.
- [15] D. Flanagan, "JavaScript: The Definitive Guide," O'Reilly Media, 7th Edition, 2020.

# Appendix

---

## A.1 Installation Instructions

### A.1.1 System Requirements

- Java 17 or higher (OpenJDK or Oracle JDK)
- Apache Maven 3.8 or higher
- Stable internet connection (for live scraping)
- 4 GB RAM minimum
- Any modern web browser (Chrome, Firefox, Edge)

### A.1.2 Installation Steps

26. Extract the project ZIP: purchaseiq\_v2\_fixed.zip
27. Open Command Prompt / Terminal
28. Navigate to backend: cd path/to/purchaseiq/backend
29. Start the server: mvn spring-boot:run
30. Wait for: Started Application in X.XXX seconds
31. Open frontend/dashboard.html in your browser
32. Click 'Scrape All Categories' to fetch live data

## A.2 Sample Dataset Format

The H2 database is automatically created at ./purchaseiq\_data.mv.db. To view it:

- Open browser: http://localhost:8080/h2-console
- JDBC URL: jdbc:h2:file:./purchaseiq\_data
- User: sa | Password: (empty)
- Click Connect

Sample SQL queries:

```
SELECT * FROM PRODUCTS WHERE RISK_LEVEL = 'CRITICAL' ORDER BY RISK_SCORE  
DESC;  
SELECT CATEGORY, AVG(RISK_SCORE) FROM PRODUCTS GROUP BY CATEGORY;  
SELECT COUNT(*) FROM PRODUCTS WHERE IN_STOCK = FALSE;
```

## A.3 API Reference

### A.3.1 BooksToScrapeScrapper Class

```
public class BooksToScrapeScrapper {
    // Scrape all 15 genre categories
    public List<Product> scrapeAllCategories();
    // Scrape a single genre with pagination
    public List<Product> scrapeGenreWithPagination(String url, String
genre);
}
```

### A.3.2 RiskAnalysisService Class

```
public class RiskAnalysisService {
    // Analyse all products in database
    public void analyzeAllProducts();
    // Analyse a single product (returns updated entity)
    public Product analyzeProduct(Product p, Map<String, CategoryStats>
stats);
}
```

## A.4 Additional Tables

### A.4.1 Complete Risk Flag Distribution

Table A.1 presents the complete distribution of risk flags triggered across 1,000 scraped books.

**Table A.1: Complete Risk Flag Distribution**

Risk Flag	Trigger Condition	Count	% of Books
OVERPRICED	Price > 2x category average	187	18.7%
LOW_RATING	Rating <= 2 stars	149	14.9%
AVERAGE_RATING	Rating = 3 stars (below acceptable threshold)	389	38.9%
OUT_OF_STOCK	Product currently unavailable	150	15.0%
SUPPLY_RISK	Category out-of-stock ratio > 30%	203	20.3%
HIGH_COMPETITION	Category contains > 50 books	142	14.2%
BELOW_CATEGORY_AVERAGE	Product rating > 1.0 below genre average	211	21.1%
SUSPICIOUSLY_CHEAP	Price < 50% of category average	43	4.3%
NO_RATING	Product has no rating data	0	0.0%

## A.5 Glossary

**CSS Selector:** A pattern-matching syntax for identifying HTML elements based on tag names, class names, IDs, and attribute values.

**Composite Risk Score:** A weighted combination of multiple independent risk dimension scores into a single normalised metric (0-100).

**DOM (Document Object Model):** The hierarchical tree representation of an HTML document, enabling programmatic traversal and querying.

**Embedded Database:** A database engine that runs within the same process as the host application, requiring no separate server process.

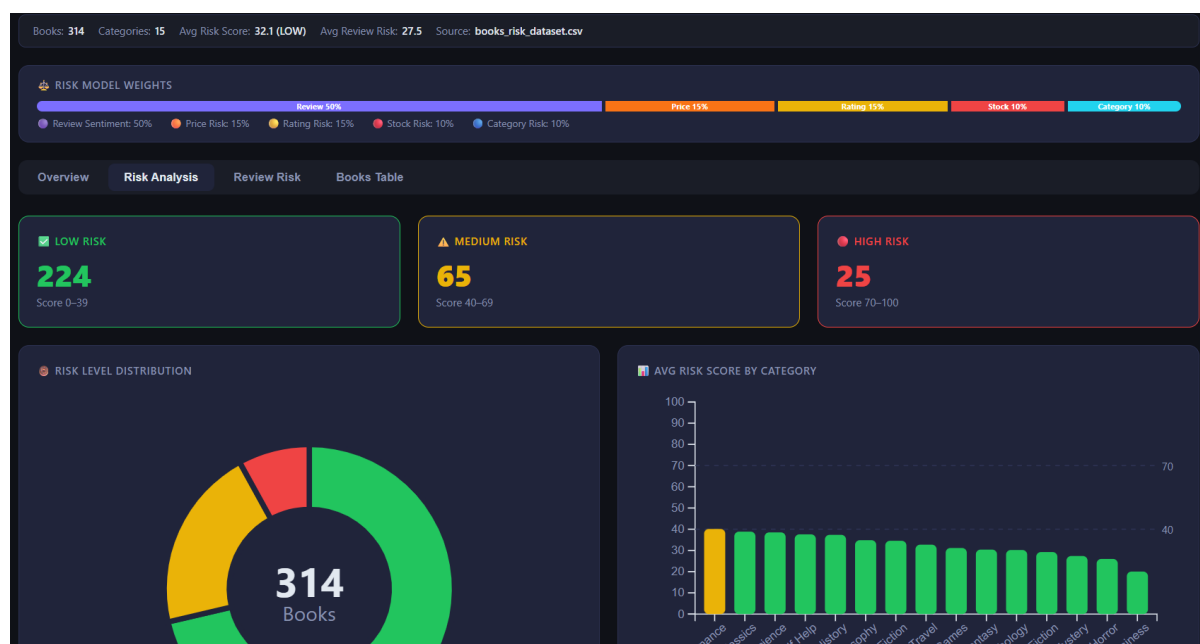
**Jsoup:** A Java HTML parsing library implementing the HTML5 specification with CSS4 selector querying support.

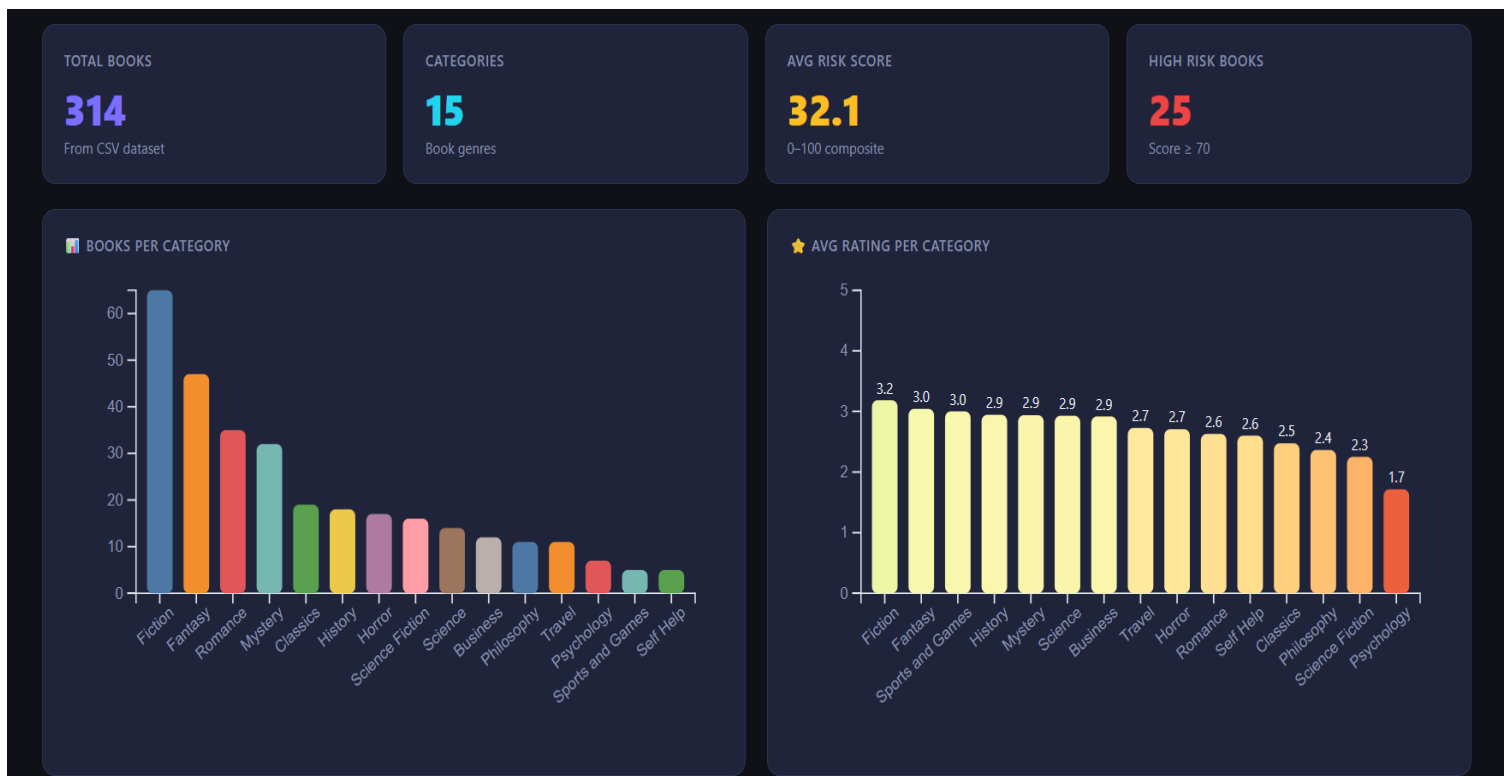
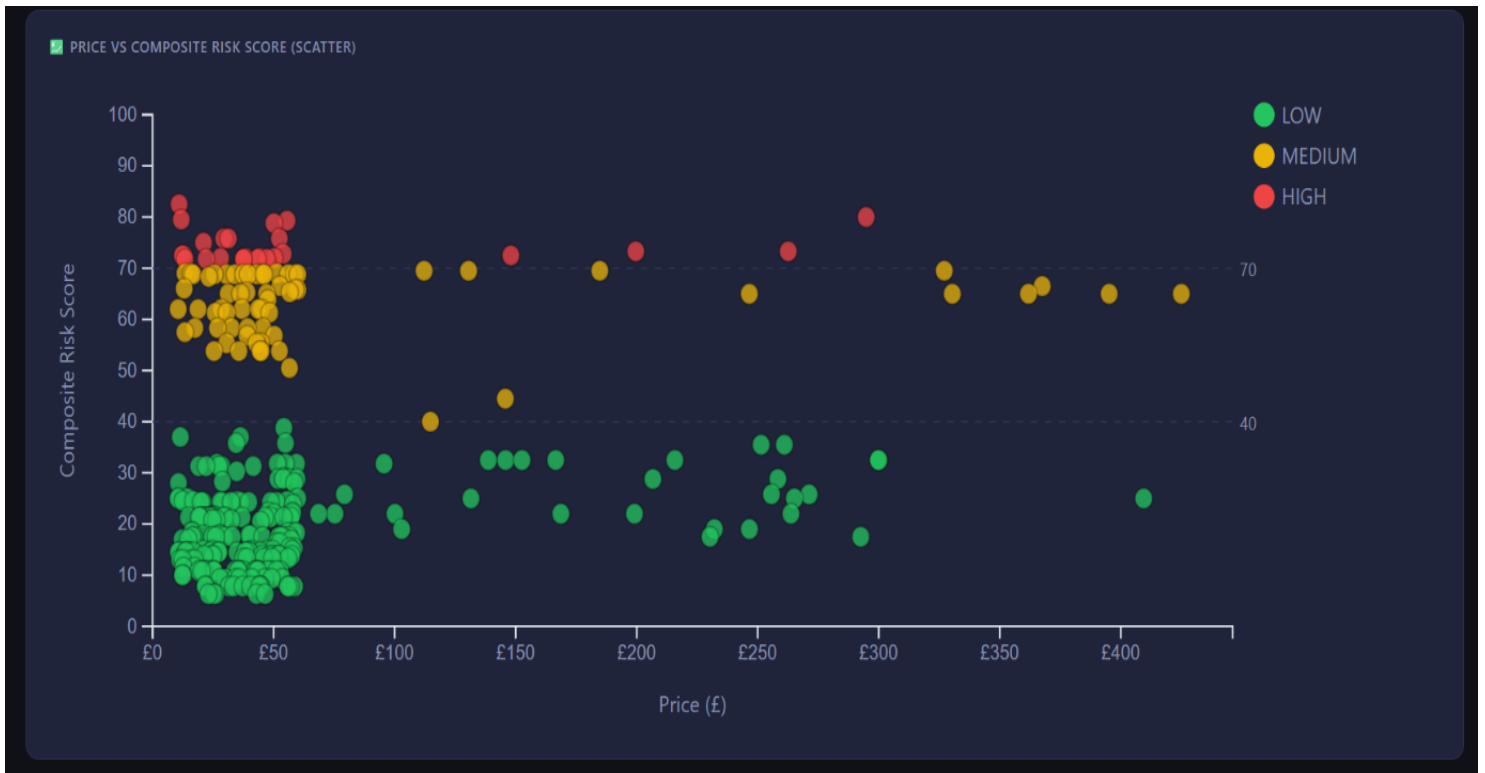
**Pagination:** The mechanism by which large datasets are split across multiple web pages, each requiring a separate HTTP request.

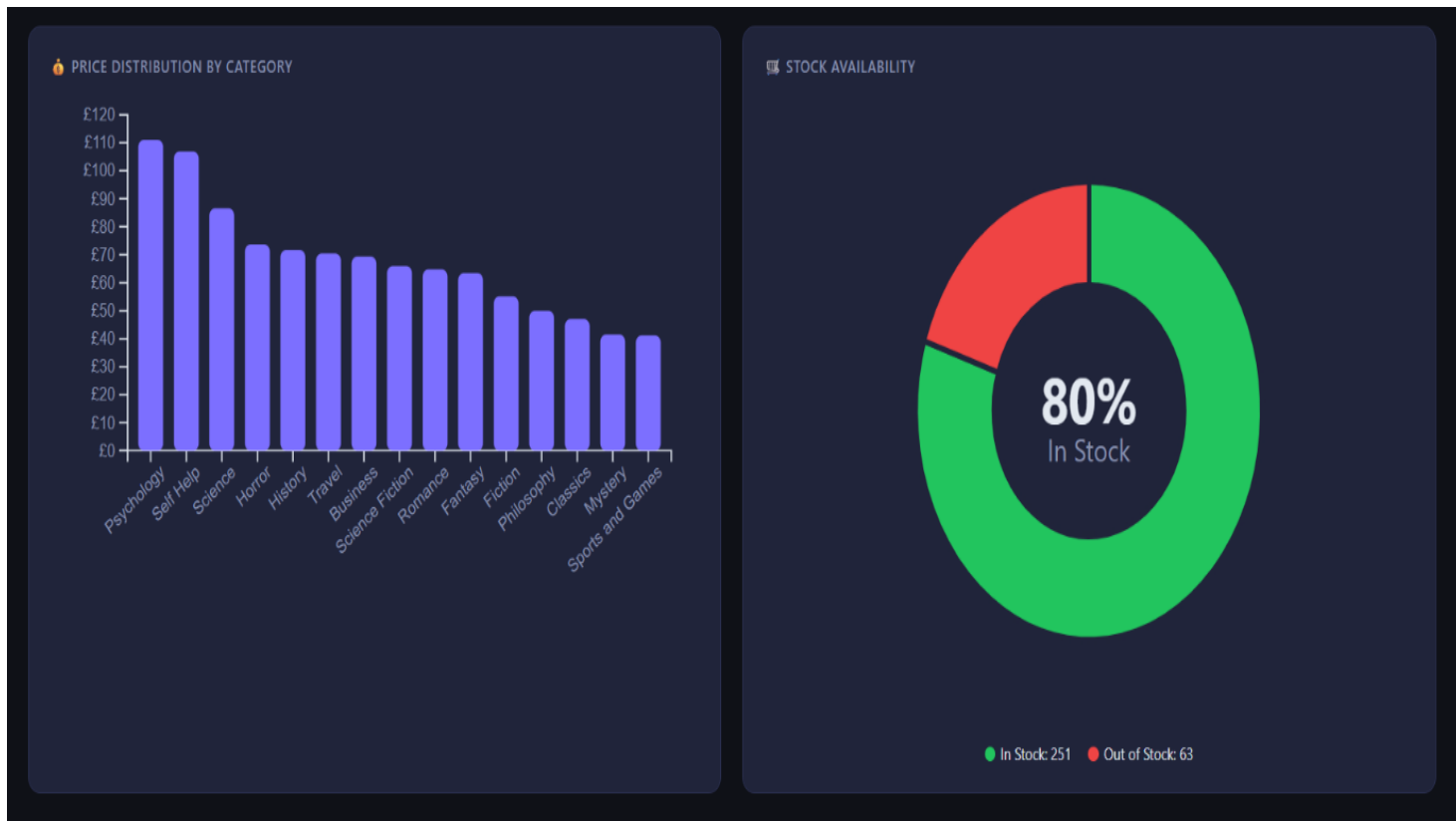
**Risk Flag:** A human-readable string describing a specific risk condition triggered for a product (e.g., OVERPRICED, LOW\_RATING).

**Spring Boot:** A Java framework that simplifies application development through auto-configuration and embedded servlet containers.

**Web Scraping:** The automated extraction of structured data from publicly accessible web pages via HTTP requests and HTML parsing.







**ALL BOOKS — REVIEW TEXT & REVIEW RISK SCORE**

Search title or review... All Sentiment Flags 314 books

TITLE	CATEGORY	REVIEW TEXT	REVIEW RISK	SENTIMENT FLAG	COMPOSITE SCORE	LEVEL
Sophie's World	Philosophy	Perfect condition, great price, and arrived before the estimated date.	0	Positive	13.0	LOW
The Death of Humanity; and the ...	Philosophy	Book was exactly as described, great condition and prompt delivery.	0	Positive	18.3	LOW
The Stranger	Philosophy	Slow delivery, took over 3 weeks to arrive.	95	Negative	58.3	MEDIUM
Proofs of God: Classical Argumen...	Philosophy	Excellent quality print, pages are crisp and clear. Very satisfied.	0	Positive	38.8	LOW
Kierkegaard: A Christian Missiona...	Philosophy	Wrong book delivered twice. Still waiting for the correct one.	95	Negative	71.8	HIGH
At The Existentialist Cafe	Philosophy	Highly satisfied! The book was well-wrapped and in mint condition.	0	Positive	9.3	LOW
Critique of Pure Reason	Philosophy	Arrived in 2 days, pristine condition. Couldn't be happier.	0	Positive	24.3	LOW
Run Spot Run: The Ethics of Keepi...	Philosophy	Excellent quality print, pages are crisp and clear. Very satisfied.	0	Positive	24.3	LOW