

# HPDP\_Project Report.pdf

*by Tan Jun Yuan*

---

**Submission date:** 13-May-2025 09:28PM (UTC-0700)

**Submission ID:** 2675452190

**File name:** HPDP\_Project\_Report.pdf (5.08M)

**Word count:** 4867

**Character count:** 30981



Department of Computer Science  
Faculty of Computing

## Optimizing High-Performance Data Processing for Large-Scale Web Crawlers

<b>Programme</b>	: Bachelor of Computer Science ( <i>Data Engineering</i> )
<b>Subject Code</b>	: SECP3133
<b>Subject Name</b>	: High Performance <b>Data</b> Processing
<b>Session-Sem</b>	: 2024/2025-2

**Prepared by** : BERNICE LIM JING XUAN (A22EC0038)  
KEK JESSLYN (A22EC0057)  
TAN JUN YUAN (A22EC0107)  
NAVACHANDER NAVASANTAR (A22EC0226)

**Section** : 01

**Lecturer** : Dr Mohd Shahizan bin Othman

**Date** : 27-04-2025

## 2 Table of Contents

<b>1.0 Introduction.....</b>	<b>1</b>
1.1 Background of the project.....	1
1.1.1 Web Scraping.....	1
1.1.2 Data Processing.....	1
1.1.3 Optimisation Process.....	2
1.2 Objectives.....	3
1.3 Target website and data to be extracted.....	3
<b>2.0 System Design &amp; Architecture.....</b>	<b>4</b>
2.1 Description of architecture.....	4
2.2 Tools and frameworks used.....	5
2.3 Roles of team members.....	5
<b>3.0 Data Collection.....</b>	<b>6</b>
3.1 Crawling method.....	6
3.2 Number of records collected.....	6
3.3 Ethical considerations.....	7
<b>4.0 Data Processing.....</b>	<b>8</b>
4.1 Cleaning methods.....	8
4.2 Data structure.....	8
4.3 Transformation and formatting.....	8
<b>5.0 Optimization Techniques.....</b>	<b>9</b>
5.1 Methods used: multithreading, multiprocessing, Spark, etc.....	9
5.2 Code overview or pseudocode of techniques applied.....	10
<b>6.0 Performance Evaluation.....</b>	<b>13</b>
6.1 Before vs after optimization.....	13
6.2 Comparison of Code Execution Time, Peak Memory Usage, CPU usage and Throughput.....	13
6.3 Charts and graphs.....	16

<b>7.0 Challenges &amp; Limitations.....</b>	<b>18</b>
7.1 What didn't go as planned.....	18
7.2 Any limitations of your solution.....	18
<b>8.0 Conclusion &amp; Future Work.....</b>	<b>19</b>
8.1 Summary of findings.....	19
8.2 What could be improved.....	19
<b>References.....</b>	<b>21</b>
<b>Appendices.....</b>	<b>21</b>
<b>Sample code snippets.....</b>	<b>22</b>
<b>Screenshots of output.....</b>	<b>22</b>
<b>Links to full code repo or dataset.....</b>	<b>30</b>

## 1.0 Introduction

### 1.1 Background of the project

In the era of big data, high-performance computing (HPC) plays a critical role in enabling the efficient processing of vast volumes of information from web sources. Web data extraction, or web scraping, has become a fundamental technique for data collection in fields such as e-commerce analysis, sentiment analysis and market research. However, handling large-scale web data introduces significant challenges, including performance bottlenecks, ethical scraping practices and managing crawl delays. To address these challenges, modern scraping systems increasingly incorporate multithreading, multiprocessing and distributed processing techniques to enhance scalability and efficiency.

This project is designed to provide students with practical, hands-on experience in large-scale web data processing using HPC principles. By designing, developing and optimising a web crawler capable of extracting at least 100,000 structured records, students gain insight into real-world technical and ethical challenges associated with web scraping. Furthermore, the project emphasises the importance of system optimisation, particularly through the comparison of different data processing frameworks, thus strengthening critical thinking skills essential for data science professionals.

#### 1.1.1 Web Scraping

This project focuses on collecting and preparing product data from Lazada Malaysia, specifically targeting women-related categories such as Beauty & Skincare, Health & Wellness, Home & Living, Home Appliances, Mother & Baby, Stationery, and Women's Fashion. The main objective is to obtain a clean and structured dataset that can later be used for further analysis or machine learning tasks.

The first step involves web scraping, where product data is automatically collected from the selected subcategories on Lazada. Each subcategory's data is then stored separately in seven Excel files for better organisation. There are a total of 115090 rows of data that have been collected. Once the data is collected, it is uploaded to Google Colab for preprocessing.

#### 1.1.2 Data Processing

In the preprocessing phase, the first task is data integration, where all seven Excel files are combined into a single dataset. To ensure consistency, all string-based fields such as product names are standardised to uppercase formatting. This helps avoid issues caused by inconsistent capitalization during analysis, such as "lotion" and "Lotion" being treated as different items.

Next, we convert important numerical fields like quantity sold and total reviews into numeric data types. This step is crucial because numeric values are required for proper data analysis, such as outlier detection and calculation needed for grouping items into categories.

After ensuring that all data is in the correct format and structure, we handle missing values. For string fields, missing values are filled with "N/A" to clearly indicate unavailable information, while missing numeric fields are filled with 0. This approach ensures that the dataset remains complete without causing errors in future computations. For example, a missing review count is more safely treated as zero than left blank.

Duplicate records are then detected and removed to avoid repetition and ensure data accuracy. Once all the cleaning steps are completed, the final, clean dataset is exported into a CSV file, ready for the optimisation process.

### 1.1.3 Optimisation Process

The second phase of the project focuses on optimising the cleaned dataset obtained from the initial preprocessing stage. The objective of this phase is to group and analyse products based on pricing tiers, popularity levels, and market performance by location in order to derive meaningful insights and support further analytical tasks.

The first optimisation step involves the categorisation of products into four pricing tiers, which are budget-friendly, affordable, mid-range and premium. Prior to grouping, all records with a price value of RM0 were removed, as such entries are considered illogical or erroneous. Outliers within the price field were then identified. Upon evaluation, these outliers were deemed plausible and were therefore retained. The minimum and maximum prices (excluding outliers) were calculated and used to define the thresholds for each pricing group. Subsequently, all products, including those with outlier prices, were assigned to the appropriate pricing category based on the established range.

The second optimization focuses on product popularity, measured by the total number of reviews. In this stage, outliers were detected and removed to minimise distortion in the analysis. The adjusted minimum and maximum review counts were then used to determine suitable group boundaries. Products were classified into four popularity levels, which are least popular, below average, above average and most popular, based on their total review count.

The final optimization step involves evaluating product performance by location. Products were grouped according to their listed locations, with relevant attributes such as product price and quantity sold. For each location, the average product price and total quantity sold were computed. These figures were used to estimate market performance, calculated by multiplying the average price by the total quantity sold. Locations were then ranked from highest to lowest based on this performance indicator, enabling identification of regions with the strongest sales activity.

## 1.2 Objectives

8

The main objectives of this project are as follows:

- To develop a web crawler capable of extracting a minimum of 100,000 structured records from a targeted Malaysian e-commerce website.
- To apply high-performance computing techniques, including multithreading, multiprocessing and distributed processing, to optimise the efficiency and scalability of the web crawling and data processing systems.
- To implement ethical web scraping practices by respecting crawl delays and website usage policies.
- To conduct a comparative performance analysis of different data processing frameworks (Pandas, Polars and PySpark) based on the time consumed during data processing.
- To enhance students' technical proficiency, critical thinking in system optimization and collaborative skills in a diverse team environment.

## 1.3 Target website and data to be extracted

For this project, Lazada Malaysia (<https://www.lazada.com.my/>) was selected as the target website. Lazada is one of the leading e-commerce platforms in Southeast Asia, offering a wide range of products across multiple categories. The focus of the data extraction is on products under the "Women" category, which includes the following subcategories: Women's Fashion, Stationery, Mother and Baby, Home and Living, Health and Wellness and Beauty and Care. The fields extracted for each product are:

- **Product Name:** The title or description of the product as displayed on the website.
- **Location:** The seller's or product's listed location.
- **Quantity Sold:** The number of units sold, indicating the popularity of the product.
- **Price:** The listed selling price of the product.
- **Total reviews:** The total number of customer ratings received by the product.

Data scraping was carried out by applying a mix of Python libraries and tools such as BeautifulSoup, Selenium, Requests for complete data extraction. The stocks of data collected were then manipulated using pandas polars and PySpark, the processing time compared to evaluate performance enhancement in varied optimisation techniques.

## 2.0 System Design & Architecture

### 2.1 Description of architecture

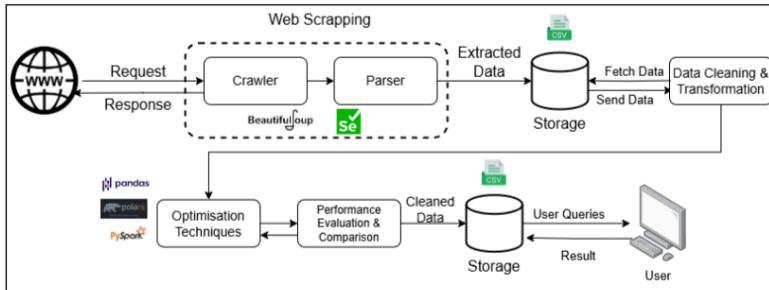


Figure 1: Web Crawler System Architecture

This project describes the design and the implementation of web crawler system for data extraction mechanism, cleaning, optimisation and analysis. In this project data under use was acquired from Lazada Malaysia with the particular interest of the Women's category. The obtained dataset is termed "Women's Purchase Analysis".

The system starts when it sends a request to the target website and get the corresponding response.

- The Crawler component traverses the web pages in an organized manner, therefore, to gather pertinent data.
- The parser, by applying libraries, such as BeautifulSoup, and Selenium, retrieves structured data from the gathered content on the web.
- Then the extracted data is saved to a CSV file, for the subsequent process.

Following data extraction:

- A data cleaning and transformation is performed in retrieving the extracted data, removing inconsistencies, missing values and standardising the dataset to maintain quality and consistency.
- Cleaned data is stored in a structured form, separately.

For performance enhancement:

- Three libraries of Pandas, Polars and PySpark are applied so as to enhance the speed in processing data.
- A performance analysis and comparison is presented to run and compare the effectiveness of these selected optimisation methods.

Last:

- Cleaned and optimised data is provided for user queries.
- The users can interact interacting with the system by sending a query and the system will generate the requested analysis results based on the processed data set.

This system guarantees full workflow in data acquisition to effective data retrieval providing full analysis of women's purchase trends on Lazada.

## 2.2 Tools and frameworks used

The following tools and framework were used during the project:

- **BeautifulSoup**: A Python library applied for parsing HTML and XML documents. It made it easier to retrieve the targeted information from the Lazada web pages as the web scraping phase.
- **Selenium**: A web automation tool which was used to communicate with dynamic pages and deal with content that needed users to scroll and click to completely load before extraction.
- **Python**: The foremost language that was used for implementation of the web scraping, data cleaning and data analysis function.
- **Pandas**: A Python library for manipulating and examining data. It was used for cleaning, transformation and processing of the extracted data effectively.
- **Polars**: A fast Rust implementation of a Dataframe library, an alternative to Pandas for Dataframe manipulation for when dealing with bigger datasets.
- **PySpark**: The Python API for Apache Spark, developed to optimize processing of larger datasets, and distributed data operations to enhance performance and scalability produces considerably fewer errors.

## 2.3 Roles of team members

Bernice Lim Jing Xuan	<ul style="list-style-type: none"><li>• Project planning and management</li><li>• Developed optimization code using pandas</li></ul>
Kek Jesslyn	<ul style="list-style-type: none"><li>• Web crawling and scraping</li><li>• Developed scrapers using BeautifulSoup and Selenium</li><li>• Conducted performance comparison testing using laptop</li></ul>
Navachander Navasantar	<ul style="list-style-type: none"><li>• Developed optimization code using PySpark</li><li>• Assisted in report documentation</li></ul>
Tan Jun Yuan	<ul style="list-style-type: none"><li>• Developed optimization code using polars</li><li>• Assisted in report work planning and documentation</li></ul>

*Table 1: Roles of Team Members*

## **3.0 Data Collection**

### **3.1 Crawling method**

The web scraping script functioned with Python tools that used Selenium to control browsers and BeautifulSoup to extract HTML data. Our scraping solution undertook multiple page operations in female-oriented Lazada product sections including fashion and skincare along with wellness and home and baby and stationery selections.

- The scraping system detected page count through pagination components before it automatically visited all accessible pages in each category.
- The scraper included rate-limiting functions that simulated human behaviors by introducing random sleep time between 2.5 and 5 seconds when triggering actions such as page loading or "next" button clicks.
- The scraper relies on async handling while using Selenium's WebDriverWait to monitor page element load times as a method to avoid incomplete content during processing. Manual detection occurs when CAPTCHAs appear since a temporary stop occurs for user confirmation.

### **3.2 Number of records collected**

A total of 115090 records were obtained from 36 distinct URLs that covered multiple pricing segments and sections including Women's Fashion, Beauty Skincare, Health Wellness, Home Living, Home Appliances, Mother & Baby and Stationery. Each record consists of:

- Product name
- Price
- Seller location
- Quantity sold
- Total reviews

### **3.3 Ethical considerations**

The data scraping operation abided by all ethical standards throughout the process.

- Research and academic needs formed the sole basis for the scraping activities.
- The company Lazada has not declared any restrictions on data scraping for our collected data so we protected their server by delaying requests while managing concurrent requests at a moderate level.
- The researchers manually dealt with CAPTCHA interruptions to prevent any security bypassing situations.
- Our data scraping operations did not involve any wishes of personalized user information.
- We would either have accessed the Lazada API or attempted to obtain permission through clear terms if the platform provided either option.

## **4.0 Data Processing**

### **4.1 Cleaning methods**

The following cleaning techniques were applied to ensure the accuracy, completeness and consistency of the raw data:

- Merged seven individual data frames into a single unified dataframe.
- Verified the combined dataset by checking the total number of rows and columns.
- Removed irrelevant text and symbols from fields such as "Quantity Sold", including handling shorthand notations (example: converting "1.3K" to 1300).
- Converted non-numeric fields into appropriate numeric types where necessary, using error coercion to handle invalid values.
- Identified and handled missing data:
  - Replaced missing numeric entries with 0.
  - Replaced missing string entries with "N/A".
- Detected and removed duplicate rows to prevent redundancy and ensure data integrity.

### **4.2 Data structure**

The final cleaned dataset was stored in the CSV (Comma-Separated Values) format, providing a lightweight and widely supported structure suitable for data processing and analysis tasks.

### **4.3 Transformation and formatting**

After cleaning, the following transformations and formatting operations were performed to prepare the data for analysis and optimisation:

- Standardised all string fields by converting text to uppercase for consistency across entries.
- Ensured numeric fields, such as "Quantity Sold" and "Number of Ratings," were properly cast to integer (Int64) data types.
- Formatted the dataset to maintain uniform data types across all columns, facilitating smoother downstream processing.
- Structured the data to eliminate inconsistencies, enabling compatibility with optimisation libraries such as pandas, polars and PySpark.

## **5.0 Optimization Techniques**

### **5.1 Methods used: multithreading, multiprocessing, Spark, etc.**

To optimize the data processing step, three python libraries were utilized; Pandas, polars and Pyspark libraries. Pandas was initially used as a benchmark because it is straightforward, and has many powerful data manipulation abilities. However, as Pandas functions in a single threaded approach it was found to have limitations with large datasets. To ensure a faster process, Polars was introduced. Polars supports multithreading and lazy evaluation which means that operations can compute across several CPU cores at once, and consequently are faster than Pandas.

Finally, for distributed processing an experimental and open source framework written in Python that runs on top of Apache Spark called PySpark was used. PySpark allows processing data at parallel across several cores or machines hence amazingly suitable for very large datasets. Its distributed architecture and native optimisation properties enabled large scale data transformation to work efficiently. Using Pandas (single-threaded), and the multithreaded variant of Polars and PySpark (distributed), the project evaluated the performance and scalability of various optimisation methods for processing data web-scraped.

## 5.2 Code overview or pseudocode of techniques applied

```
Part 1 Data Processing and Cleaning

Prepared by : BERNICE LIM JING XUAN (A22EC0038)

Step 1: Install and Import Libraries

[ ] pip install pandas
import pandas as pd
import time
import tracemalloc
import psutil
import os

>Show hidden output

Step 2: Upload Excel files

[ ] from google.colab import files
uploaded = files.upload()

>Show hidden output

Step 3: Load Excel Files into PANDAS DataFrames

❶ flist_pandas = [pd.read_excel(file) for file in uploaded.keys()]
print("Total Dataframes in flist_pandas: ", len(flist_pandas))

>Show hidden output

Step 4: Load and Display Dataset, Checking on Total Numbers of Rows and Columns

❷ %%time
tracemalloc.start()
start_time = time.perf_counter()
total_rows = 0

flist_pandas_with_category = []
for filename, df in zip(uploaded.keys(), flist_pandas):
    match = re.search("(.*).xlsx", filename)
    category = match.group(1) if match else "Unknown"
    df["Category"] = category
    flist_pandas_with_category.append(df)

print("Total Dataframes: ", len(flist_pandas_with_category))

for df in flist_pandas_with_category:
    total_rows += df.shape[0]
    display(df.head(10))
    print("Total rows: ", df.shape[0])
    print("Total columns: ", df.shape[1])\n\n

current, peak = tracemalloc.get_traced_memory()
end_time = time.perf_counter()
tracemalloc.stop()

execution_time = end_time - start_time
throughput = total_rows / execution_time

print("----- performance -----")
print("Total rows processed: ", total_rows)
print("Code Execution time: ", execution_time, "seconds")
print("Throughput: ", total_rows / execution_time, "rows/second")
print("Current memory usage: (current / 10^6), MB")
print("Peak memory usage: (peak / 10^6), MB")

cpu_usage = psutil.cpu_percent(interval=1)
print("CPU usage: ", cpu_usage)

print("-----")
print("Total time for this cell(including time to display the performance):")
```

Figure 2: Code Overview of Pandas

```

✓ Part1 Data Processing and Cleaning
Prepared by: TAN JUN YUAN (A22EC0107)

Step 1: Install and Import Libraries
[ ] pip install polars
import polars as pl
import pandas as pd
import re
import time
import tracemalloc
import tracemalloc

>Show hidden output

Step 2: Upload Excel Files
[ ] from google.colab import files
uploaded = files.upload()
>Show hidden output

Step 3 : Load Excel Files into PANDAS DataFrames and Check Total Files being Loaded
[ ] filelist_pandas = [pd.read_excel(file) for file in uploaded.keys()]
print(f"Total File: {len(filelist_pandas)}")
>Show hidden output

Step 4 : Load and Display Dataset, Checking on Total Numbers of Rows and Columns
[ ] %%time
tracemalloc.start()
start_time = time.perf_counter()
total_rows = 0

filelist_polars = []
for filename, df in zip(uploaded.keys(), filelist_pandas):
    match = re.search(r'^(.*\.(?!\.parquet))$', filename)
    category = match.group(1) if match else "Unknown"
    pl_df = pl.from_pandas(df).with_columns(pl.lit(category).alias("category"))
    filelist_polars.append(pl_df)

print(f"Total Dataframes: {len(filelist_polars)}")

for df in filelist_polars:
    total_rows += df.shape[0]
    display(df.head(10))
    print(f"Total columns: {df.shape[1]}\n\n")
print(f"Total columns: {total_rows}\n\n")

current_peak = tracemalloc.get_traced_memory()
end_time = time.perf_counter()
tracemalloc.stop()

execution_time = end_time - start_time
throughput = total_rows / execution_time

print("----- Performance -----")
print(f"Total rows processed: {total_rows}")
print(f"Code Execution time: {execution_time:.4f} seconds")
print(f"Throughput: {throughput:.2f} rows per second")
print(f"Peak memory usage: {current_peak / 10**6:.2f} MB")
print(f"Peak memory usage: {peak / 10**6:.2f} MB")

cpu_usage = psutil.cpu_percent(interval=1)
print(f"CPU usage: {cpu_usage}%")

print("-----")
print("Total time for this cell(including time to display the performance):")

```

**Figure 3: Code Overview of Polars**

Part1 Data Processing and Cleaning

Prepared by:

Step 1: Install and Import Libraries

```
[ ] from pyspark.sql import SparkSession
from pyspark.sql.types import StringType, NumericType
from pyspark.sql.functions import col, upper, regexp_replace, when, isnan, count, lit, mean, sum as spark_sum, avg
from functools import reduce
import pandas as pd
import time
import tracemalloc
import psutil
import os
import shutil
spark = SparkSession.builder.appName("ExcelProcessing").getOrCreate()
```

Step 2: Upload Excel Files

```
[ ] from google.colab import files
uploaded = files.upload()
Show hidden output
```

Step 3 : Load Excel Files into PANDAS DataFrames and Check Total Files being Loaded

```
[ ] list_pandas = [pd.read_excel(file) for file in uploaded.keys()]
print("Total Dataframes in list_pandas: " + str(len(list_pandas)))
Show hidden output
```

Step 4 : Load and Display Dataset, Checking on Total Numbers of Rows and Columns

```
XXTime
tracemalloc.start()
start_time = time.perf_counter()
total_rows = 0

placeholders = ["\\n\\a", "\\n", "null", "", "...", "...", "\\n\\a", "\\n\\a", "\\n\\a"]
list_spark = []

for filename, df in zip(uploaded.keys(), [pd.read_excel(file) for file in uploaded.keys()]):
    match = re.search("(?i)([\\w\\.]+\\.(\\w+))", filename)
    category = match.group(1).strip() if match else "Unknown"
    df_cleaned = df.replace(placeholders, None)
    spark_df = spark.createDataFrame(df_cleaned)

    spark_df = spark_df.select([
        when(col(c).isin(placeholders), None).otherwise(col(c)).alias(c)
        if spark_df.schema[c].datatype.simpleString() == "string" else col(c)
        for c in spark_df.columns
    ])

    spark_df = spark_df.withColumn("Category", lit(category))
    list_spark.append(spark_df)

print("Total Dataframes in list_spark: " + str(len(list_spark)))

for df in list_spark:
    display(df.limit(10).toPandas())
    row_count = df.count()
    total_rows += row_count
    print("Total rows: " + str(row_count))
    print("Total columns: " + str(col_count))

current_peak = tracemalloc.get_traced_memory()
end_time = time.perf_counter()
tracemalloc.stop()

execution_time = end_time - start_time
throughput = total_rows / execution_time

print("----- Performance -----")
print("Total rows processed: " + str(total_rows))
print("Execution time: (" + str(execution_time) + ".af) seconds")
print("Throughput: (" + str(throughput) + ".af) rows per second")
print("Current memory usage: (" + str(current_peak / 10**6) + ".af) MB")
print("Peak memory usage: (" + str(peak / 10**6) + ".af) MB")

cpu_usage = psutil.cpu_percent(interval=1)
print("CPU usage: " + str(cpu_usage) + "%")

print("Total time for this cell (Including time to display the performance):")
```

Show hidden output

Figure 4: Code Overview of PySpark

## 6.0 Performance Evaluation

### 6.1 Before vs after optimization

Initially, all data processing tasks were carried out using Pandas. While Pandas performed well for small datasets, it became noticeably slower as the data size increased, especially with operations like filtering, aggregation and joining. This led to delays that affected the overall workflow efficiency.

To improve performance, two faster alternatives were introduced, which are Polars and PySpark. Polars with its multi-threaded Rust backend, offered significant speed improvements for in-memory operations, while PySpark provided better handling for larger datasets through distributed processing even in local mode. After optimisation, execution times were greatly reduced, with Polars offering the fastest performance and PySpark also outperforming Pandas, although with a slight overhead due to its distributed nature. Overall, the optimisation resulted in a much faster and more scalable data processing pipeline.

### 6.2 Comparison of Code Execution Time, Peak Memory Usage, CPU usage and Throughput

Operation	Aspects	Comparisons		
		Pandas	Polars	Pyspark
Dataset Loading and Display	Code Execution Time (s)	0.6728	0.3893	113.0075
	Peak Memory Usage (MB)	1.4805	2.9626	34.7047
	CPU Usage (%)	2.5	2.0	84.8
	Throughput (rows/s)	171056.76	295636.96	1018.43
Dataset Integration	Code Execution Time (s)	0.0907	0.0207	4.8765
	Peak Memory Usage (MB)	5.6588	0.0419	0.1984
	CPU Usage (%)	3.5	2.5	12.6
	Throughput (rows/s)	1269298.78	5564845.61	23601.16
Standardization of String Data	Code Execution Time (s)	0.6550	0.0447	4.2237
	Peak Memory Usage (MB)	50.4867	0.0408	0.1824
	CPU Usage (%)	3.5	2.0	61.0
	Throughput (rows/s)	175721.43	2574996.64	27248.92
Convert “Total Reviews” to correct data type	Code Execution Time (s)	2.8809	0.0260	5.1355
	Peak Memory Usage (MB)	12.7351	0.0420	0.1679

(int)	CPU Usage (%)	4.0	2.5	21.6
	Throughput (rows/s)	39949.87	4434821.69	22410.57
Convert “Quantity Sold” to correct data type (int)	Code Execution Time (s)	3.9448	0.0843	4.1328
	Peak Memory Usage (MB)	19.3337	0.3887	0.1938
	CPU Usage (%)	25.6	2.0	6.0
	Throughput (rows/s)	29174.87	1365230.27	27847.80
Check and Handle Missing Values	Code Execution Time (s)	0.3365	0.0300	18.2981
	Peak Memory Usage (MB)	19.8926	0.0498	0.2372
	CPU Usage (%)	3.5	2.0	12.5
	Throughput (rows/s)	341984.53	3836148.05	6289.71
Check and Handle Duplicates	Code Execution Time (s)	0.6107	0.0859	65.1819
	Peak Memory Usage (MB)	20.1169	0.0745	2.3469
	CPU Usage (%)	3.5	2.5	13.6
	Throughput (rows/s)	188440.81	1339841.70	1742.75
Export cleaned dataset file for optimization	Code Execution Time (s)	4.50939	0.0592	10.6576
	Peak Memory Usage (MB)	3.8787	0.0953	0.0880
	CPU Usage (%)	3.5	6.0	75.1
	Throughput (rows/s)	25191.72	1918315.45	10658.73

Table 2: Comparison between Data Processing and Cleaning Techniques

Operation	Aspects	Comparisons		
		Pandas	Polars	Pyspark
Grouping Products into 4 categories based on price (Budget Friendly, Affordable, Mid-Range and Premium Price)	Code Execution Time (s)	0.6036	0.5189	22.6163
	Peak Memory Usage (MB)	19.2549	0.1225	0.3516
	CPU Usage (%)	88.6	90.0	5.5
	Throughput (rows/s)	188187.82	218918.57	5022.74
Grouping Products into 4 categories based on ‘Total Reviews’ (Least, Below Average, Above Average and Most Popular)	Code Execution Time (s)	0.9005	0.2981	13.5979
	Peak Memory Usage (MB)	15.0830	0.0664	0.3121
	CPU Usage (%)	86.9	100.0	19.0
	Throughput (rows/s)	126141.24	381122.98	8353.95
Evaluate and Rank Market Performance based on “Quantity Sold” for each “Location”	Code Execution Time (s)	2.1653	0.1121	1.8099
	Peak Memory Usage (MB)	6.2297	0.0222	0.1632
	CPU Usage (%)	95.5	100.0	56.9
	Throughput (rows/s)	52461.36	1013159.21	62763.71

Table 3: Comparison between Data Optimization Techniques

#### Conclusion:

Overall performance between the libraries: **Polars > Pandas > PySpark**

- **Polars:** Polars delivered superior performance than PySpark on big datasets thanks to its Rust-based foundation coupled with threading capabilities and columnar data structure.
- **Pandas:** The system efficiency of Pandas was high but it utilised more system memory and took longer to process data.
- **PySpark:** PySpark registered the slowest performance since it required high resource allocation latency combined with substantial initialisation overhead.

### 6.3 Charts and graphs

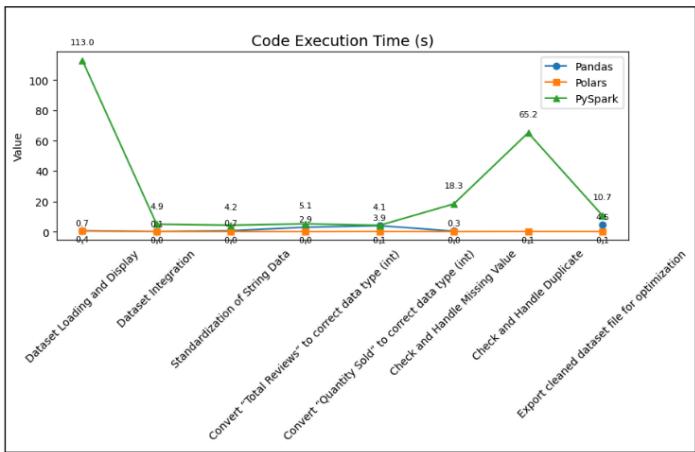


Figure 5: Line graph for Code Execution Time (s)

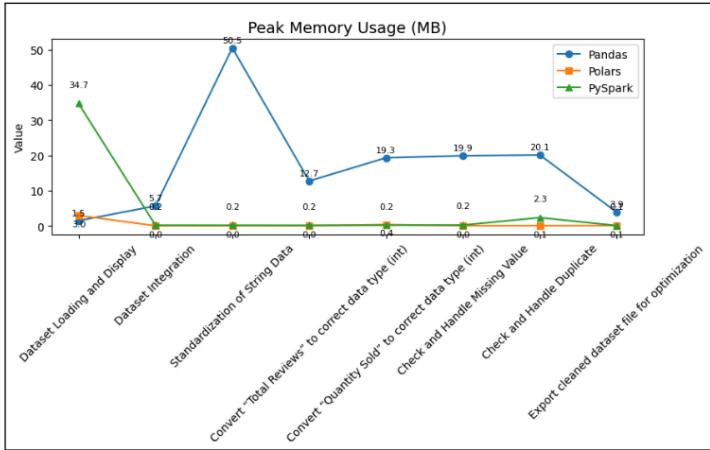


Figure 6: Line graph for Peak Memory Usage (MB)

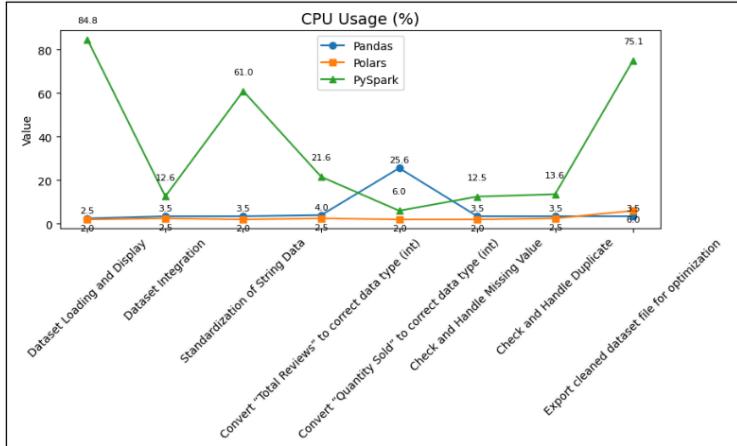


Figure 7: Line graph for CPU Usage (%)

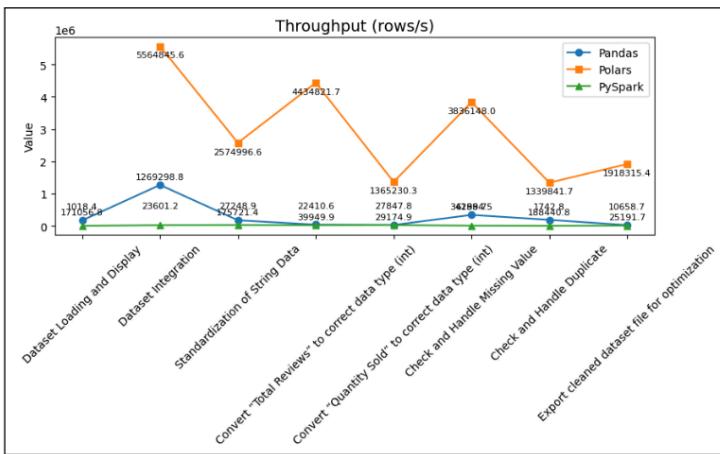


Figure 8: Line graph for Throughput (rows/s)

## **7.0 Challenges & Limitations**

### **7.1 What didn't go as planned**

Initially, the project intended to use only Requests and BeautifulSoup for web scraping. However, Lazada's dynamic JavaScript-driven content made it impossible to retrieve the necessary information with these tools alone. Therefore, Selenium was incorporated to handle the dynamic loading of data. However, using Selenium brought new issues, particularly the need for ChromeDriver installation. Since Google Colab could not successfully run ChromeDriver, the team had to shift the environment to Visual Studio Code on local machines. This transition increased the setup time and complexity.

Another major challenge was Lazada's frequent CAPTCHA verifications. The team had to manually solve CAPTCHAs during the scraping process, which greatly slowed down data collection. Scraping had to be done during weekends, requiring long hours in front of the computer to monitor and complete the verification steps.

During the data processing stage, differences in library capabilities also posed problems. For example, Polars could not directly read Excel files. To resolve this, the data was first read into Pandas and then converted into a Polars DataFrame, adding extra steps and minor inefficiencies to the workflow.

### **7.2 Any limitations of your solution**

Despite overcoming many challenges, some limitations remained. The most significant was the reliance on manual CAPTCHA handling, which prevented full automation and reduced scraping efficiency. Running Selenium locally on Visual Studio Code also restricted collaboration, as each member needed to configure their own environment separately. This limited the flexibility that cloud platforms like Colab would have provided.

In addition, extra conversion steps between data formats affected the purity of the performance comparisons between Pandas, Polars and PySpark. Although functional, it introduced slight variations in the benchmarking results.

Finally, the need for manual oversight and longer scraping times meant that the project could not easily scale to even larger datasets or implement more advanced scraping strategies like headless browsing or automated CAPTCHA solving.

## **8.0 Conclusion & Future Work**

### **8.1 Summary of findings**

This project sought to improve high-efficiency data processing for massive web crawling of Lazada Malaysia's Women category. Our group collected over 115,000 product listings by utilising the tools named Selenium and BeautifulSoup to navigate the site and pull out information such as product titles, price, number sold, seller locations, and customer reviews. We filtered and arranged the data once it was gathered, prepared for analysis.

We then contrasted and examined three distinct libraries to determine which one was most appropriate to process large-scale data effectively:

- Pandas, though easy to use and effective for small to midsize datasets, showed worsening performance and more memory usage with our large dataset and thus turned out to be less effective for large-scale data processing tasks.
- Polars, with a Rust backend and lazy evaluation approach, exhibited excellent performance gains in speed and memory consumption. It offered a perfect harmony between usability and functionality for moderately sized datasets.
- PySpark was typified by scalability, distributed processing, and appropriateness for large-scale data projects. Still, its configurational complexity and associated overhead rendered it less suitable for small or medium-size projects.

Overall, this project provided our group with hands-on practice, from data collections and cleaning to performance testing of modern data processing tools. Most significantly, we learned how to assess and choose suitable technologies depending on certain requirements such as dataset size, processing time, and memory consumption. This project not only improved our technical expertise in Python and high performance libraries but also improved our critical thinking, collaboration, and problem solving skills, equipping us for upcoming projects in data engineering and data analysis.

## 8.2 What could be improved

While the current implementation successfully achieved its primary goals, there are several meaningful ways the project can be expanded and enhanced in the future:

1. Automate CAPTCHA Handling for seamless handling

Manual CAPTCHA solving during scraping delays the data collection process. 2Captcha browser extension allows skipping reCAPTCHAs. This extension automates the process of solving reCAPTCHA, making it easier and faster for users to bypass these verifications [1]. Incorporating automated services such as 2Captcha can automate the process with machine learning powered CAPTCHA solvers.

2. Make the most of GPU Powered Libraries for better performance

Libraries such as RAPIDS cuDF by NVIDIA can leverage GPUs to greatly accelerate data processing operations. For instance, cuDF accelerates pandas with zero code changes and brings greatly improved performance [2]. This is particularly useful when working with large numeric datasets.

3. Use machine learning for further analysis

The structured data can be utilised to implement machine learning algorithms for customer segmentation. Machine learning methodologies are a great tool for analyzing customer data and finding insights and patterns. Artificially intelligent models are powerful tools for decision-makers. They can precisely identify customer segments, which is much harder to do manually or with conventional analytical methods [3]. For instance, K-Means is efficient machine learning when it comes to solving data cluster problems.

By adopting these improvements, the project has the potential to evolve into a robust, scalable and intelligent data pipeline capable of supporting practical, data-driven decisions in the e-commerce landscape.

## References

- [1] Captcha Solver: reCAPTCHA solver and captcha solving service. Bypass captchas using the best auto captcha solver online API - 2Captcha. (2025). 2captcha.com.  
<https://2captcha.com/>
- [2] Open GPU Data Science. (n.d.). RAPIDS. <https://rapids.ai/>
- [3] Kumar, D. (2021, June 18). Implementing Customer Segmentation Using Machine Learning [Beginners Guide]. Neptune.ai.  
<https://neptune.ai/blog/customer-segmentation-using-machine-learning>

## Appendices



Figure 9: Group photo

## Sample code snippets

### 1. Web Scraping: [Codes](#)

### 2. Pandas:

- [Part 1](#)
- [Part 2](#)

### 3. Polars:

- [Part 1](#)
- [Part 2](#)

### 4. PySpark:

- [Part 1](#)
- [Part 2](#)

### 5. Graph: [Codes](#)

## Screenshots of output

Product Name	Price	Location	Quantity Sold	Total Reviews	Category
Classic Gel Pen 1008 Black Blue Red 3x0.5mm ...	0.39	Penang	16.2K sold	(482)	Stationery
Sensis Eraser Stationery Non-Dandruff Eraser C...	0.96	Selangor	1.1K sold	(103)	Stationery
Total rows: 22316 Total customers: 8					
<hr/>					
Upsee Women's Fashion Heat Resistant Long Cut...	17.13	China	763 sold	(181)	Women's Fashion
Summer women's comfortable high-end lace dress ...	13.00	China	54 sold	(14)	Women's Fashion
Summer New Sports Suit Women's Fashion Slim P...	22.75	China	286 sold	(62)	Women's Fashion
Daring Backless V-Neck Dress Short Mini Skirt ...	19.72	Natl.	270 sold	(37)	Women's Fashion
Chi Luu Birthday Dress Christmas Lounge Series	12.18	Natl.	486 sold	(38)	Women's Fashion
LOMOCC Women's Fashion Dress + Outer Cardigan...	15.23	China	1.4K sold	(455)	Women's Fashion
Summer 2024 Women's Knitted Long Halter Outfit ...	22.41	Natl.	322 sold	(28)	Women's Fashion
Hot Princess Dress Set Sexy Styling Wear Red ...	17.39	Natl.	92 sold	(8)	Women's Fashion
LOMOCC Women's Fashion Dress + Outer Cardigan...	14.17	Natl.	122 sold	(42)	Women's Fashion
Adult Latin Dance Skirt Square Dance Costume H...	20.85	Natl.	382 sold	(58)	Women's Fashion
Total rows: 9509 Total customers: 6					
<hr/> <b>Performance</b> <hr/>					
Total rows processed: 115950					
Rows per second: 1,939 rows per second					
Throughput: 1703294.79 rows per second					
Current memory usage: 5.3819 MB					
Peak memory usage: 5.4001 MB					
CPU usage: 3.3%					
<hr/> Total time for this cell (including time to display the performance): 272.11 ms, avg: 87.9 ms, stdev: 7.88 ms, total: 95.7 ms Wall time: 1.01 s					

Product Name	Price	Location	Quantity Sold	Total Reviews	Category
YSAKA GLOWING FOUNDATION SPF 50+ / PA++ Fath... W...	3.99	Penang	55 sold	(9)	Beauty & Skincare
Bio-Essence Bio-Gold 24K RadianceWhiteningB...	3.50	Johor	46 sold	(10)	Beauty & Skincare
L'occitane Immortelle Divine Firming Cleansing...	1.50	Selangor	13 sold	(1)	Beauty & Skincare
YOUBUY Freude Cream Effectively Revitalize Melone...	5.45	China	13 sold	(4)	Beauty & Skincare
DT3 Jerjen Joyen Nicotinide Ampoule Facial ...	6.15	Malaka	175 sold	(42)	Beauty & Skincare
W88incore W88incore Jomtan Joytan Nicotinide Ampou...	6.19	Malaka	29 sold	(2)	Beauty & Skincare
DT3 VEZZE Men's Volcanic Mud Facial Cleanser B...	5.44	Malaka	123 sold	(48)	Beauty & Skincare
NEVER MEN ROLL ON 500ML 5.00	5.00	Wp Kuala Lumpur	25 sold	(4)	Beauty & Skincare
W88incore VEZZE Men's Volcanic Mud Facial Clea...	5.43	Malaka	41 sold	(11)	Beauty & Skincare
DT3 HYDRO Facial Cleanser Cream Whitening M...	1.93	Malaka	1.3K sold	(248)	Beauty & Skincare
Total rows: 115950 Total customers: 6					
<hr/> <b>Performance</b> <hr/>					
Total rows processed: 115950					
Rows per second: 1,939 rows per second					
Throughput: 1703294.79 rows per second					
Current memory usage: 5.3819 MB					
Peak memory usage: 5.4001 MB					
CPU usage: 3.3%					
<hr/> Total time for this cell (including time to display the performance): 272.11 ms, avg: 87.9 ms, stdev: 7.88 ms, total: 95.7 ms Wall time: 1.01 s					

Figure 10: Output Screenshot of Dataset Integration by using Pandas

Figure 9: Output Screenshot of Dataset Loading and Display by using Pandas

Product Name	Price	Location	Quantity Sold	Total Reviews	Category
SAKA GLOWING FOUNDATION SPF 50+ / FW FATH W...	3.99	PENANG	55 SOLD	0 BEAUTY & SKINCARE	
BIO-ESSENCE BIO-GOLD 24K RADIANCE WHITENING...	3.90	JOHOR	46 SOLD	10 BEAUTY & SKINCARE	
L'OCOTINNE IMMORTELLE DIVINE FOAMING CLEANSING...	1.80	SELANGOR	13 SOLD	1 BEAUTY & SKINCARE	
YOUNI FREDOLIC CREAM EFFECTIVELY REMOVE MELASMA...	5.45	CHINA	13 SOLD	0 BEAUTY & SKINCARE	
DT37 JONATAM JULYAH NEOTIMIDE AMINO ACID FACIA...	6.15	MELAKA	175 SOLD	40 BEAUTY & SKINCARE	
SH HOME VIZZI MENS VOLCANO MUD FACIAL CLEANSER...	6.19	MELAKA	29 SOLD	2 BEAUTY & SKINCARE	
DT37 VEZE MENS VOLCANO MUD FACIAL CLEANSER B...	5.44	MELAKA	129 SOLD	48 BEAUTY & SKINCARE	
NEVA MEN ROLL ON 50ML	5.00	WP KUALA LUMPUR	25 SOLD	4 BEAUTY & SKINCARE	
SH HOME VIZZI MENS VOLCANO MUD FACIAL CLEA...	5.45	MELAKA	41 SOLD	11 BEAUTY & SKINCARE	
DT37 HYNEY'S FACIAL CLEANSER CREAM WHITENING M...	1.53	MELAKA	13K SOLD	240 BEAUTY & SKINCARE	

Figure 11: Output Screenshot of Standardization of String Data by using Pandas

Product Name	Price	Location	Quantity Sold	Total Reviews	Category
SAKA GLOWING FOUNDATION SPF 50+ / FW FATH W...	3.99	PENANG	55 SOLD	0 BEAUTY & SKINCARE	
BIO-ESSENCE BIO-GOLD 24K RADIANCE WHITENING...	3.90	JOHOR	46 SOLD	10 BEAUTY & SKINCARE	
L'OCOTINNE IMMORTELLE DIVINE FOAMING CLEANSING...	1.80	SELANGOR	13 SOLD	1 BEAUTY & SKINCARE	
YOUNI FREDOLIC CREAM EFFECTIVELY REMOVE MELASMA...	5.45	CHINA	13 SOLD	0 BEAUTY & SKINCARE	
DT37 JONATAM JULYAH NEOTIMIDE AMINO ACID FACIA...	6.15	MELAKA	175 SOLD	40 BEAUTY & SKINCARE	
SH HOME VIZZI MENS VOLCANO MUD FACIAL CLEANSER...	6.19	MELAKA	29 SOLD	2 BEAUTY & SKINCARE	
DT37 VEZE MENS VOLCANO MUD FACIAL CLEANSER B...	5.44	MELAKA	129 SOLD	48 BEAUTY & SKINCARE	
NEVA MEN ROLL ON 50ML	5.00	WP KUALA LUMPUR	25 SOLD	4 BEAUTY & SKINCARE	
SH HOME VIZZI MENS VOLCANO MUD FACIAL CLEA...	5.45	MELAKA	41 SOLD	11 BEAUTY & SKINCARE	
DT37 HYNEY'S FACIAL CLEANSER CREAM WHITENING M...	1.53	MELAKA	13K SOLD	240 BEAUTY & SKINCARE	

Figure 12: Output Screenshot of Converting “Total Reviews” to int Data Type by using Pandas

Product Name	Price	Location	Quantity Sold	Total Reviews	Category
SAKA GLOWING FOUNDATION SPF 50+ / FW FATH W...	3.99	PENANG	55	0 BEAUTY & SKINCARE	
BIO-ESSENCE BIO-GOLD 24K RADIANCE WHITENING...	3.90	JOHOR	46	10 BEAUTY & SKINCARE	
L'OCOTINNE IMMORTELLE DIVINE FOAMING CLEANSING...	1.80	SELANGOR	13	1 BEAUTY & SKINCARE	
YOUNI FREDOLIC CREAM EFFECTIVELY REMOVE MELASMA...	5.45	CHINA	13	0 BEAUTY & SKINCARE	
DT37 JONATAM JULYAH NEOTIMIDE AMINO ACID FACIA...	6.15	MELAKA	175	40 BEAUTY & SKINCARE	
SH HOME VIZZI MENS VOLCANO MUD FACIAL CLEANSER...	6.19	MELAKA	29	2 BEAUTY & SKINCARE	
DT37 VEZE MENS VOLCANO MUD FACIAL CLEANSER B...	5.44	MELAKA	129	48 BEAUTY & SKINCARE	
NEVA MEN ROLL ON 50ML	5.00	WP KUALA LUMPUR	25	4 BEAUTY & SKINCARE	
SH HOME VIZZI MENS VOLCANO MUD FACIAL CLEA...	5.45	MELAKA	41	11 BEAUTY & SKINCARE	
DT37 HYNEY'S FACIAL CLEANSER CREAM WHITENING M...	1.53	MELAKA	1300	240 BEAUTY & SKINCARE	

Figure 13: Output Screenshot of Converting “Quantity Sold” to int Data Type by using Pandas

Product Name	Price	Location	Quantity Sold	Total Reviews	Category
SAKA GLOWING FOUNDATION SPF 50+ / FW FATH W...	3.99	PENANG	55	0 BEAUTY & SKINCARE	
BIO-ESSENCE BIO-GOLD 24K RADIANCE WHITENING...	3.90	JOHOR	46	10 BEAUTY & SKINCARE	
L'OCOTINNE IMMORTELLE DIVINE FOAMING CLEANSING...	1.80	SELANGOR	13	1 BEAUTY & SKINCARE	
YOUNI FREDOLIC CREAM EFFECTIVELY REMOVE MELASMA...	5.45	CHINA	13	0 BEAUTY & SKINCARE	
DT37 JONATAM JULYAH NEOTIMIDE AMINO ACID FACIA...	6.15	MELAKA	175	40 BEAUTY & SKINCARE	
SH HOME VIZZI MENS VOLCANO MUD FACIAL CLEANSER...	6.19	MELAKA	29	2 BEAUTY & SKINCARE	
DT37 VEZE MENS VOLCANO MUD FACIAL CLEANSER B...	5.44	MELAKA	129	48 BEAUTY & SKINCARE	
NEVA MEN ROLL ON 50ML	5.00	WP KUALA LUMPUR	25	4 BEAUTY & SKINCARE	
SH HOME VIZZI MENS VOLCANO MUD FACIAL CLEA...	5.45	MELAKA	41	11 BEAUTY & SKINCARE	
DT37 HYNEY'S FACIAL CLEANSER CREAM WHITENING M...	1.53	MELAKA	1300	240 BEAUTY & SKINCARE	

Figure 14: Output Screenshot of Checking and Handling Missing Values by using Pandas

**Figure 16: Output Screenshot of Exporting Cleaned Dataset File by using Pandas**

**Figure 15: Output Screenshot of Checking and Handling Duplicates by using Pandas**

Group	[Budget Friendly Price]:	Product Name	Price	Location	Quantity Sold	Total Reviews	Category
Group 1	[Budget Friendly Price]: 48982 products						
Group 1	[Affordable Price]: 10999 products						
Group 3	[Mid-Range Price]: 38837 products						
Group 4	[Premium Price]: 14021 products						
Group 4	[Budget Friendly Price]:						
54773	LIVE TRACKING#FOR SALES	0.05	KELANTAN	0	1	HOME & LIVING	
55102	EMARLUSIA 3PM HANDBAG	0.05	SELANGOR	0	9	HOME APPLIANCES	
51887	1PC 0.3MM BALLOON IN PEN MATA JANAH NEEDLE	0.10	SELANGOR	33500	1554	STATIONERY	
27877	BEIRUO THAI HERB ROLLING MASSAGE LINE, CONTACT SEAT	0.10	CHINA	10	3	HEALTH & WELLNESS	
31656	BEEFES BEEFES PLUS VITAMIN E BEAN BOOSTER	0.10	JOHOR	0	0	HEALTH & WELLNESS	
51640	SPOT BEEFES PLUS EYE CARE BOOSTER	0.10	JOHOR	0	0	HEALTH & WELLNESS	
55448	READY STOCK MICE MASCOT CASTOM MASKET	0.11	WP KUALA LUMPUR	49	1	HOME & LIVING	
55449	HAR HAIR HAWAN GANTING KARV YAKU PERASANG	0.12	WP KUALA LUMPUR	24	1	HOME & LIVING	
20205	PENUTUP BOTOL CAP ARTMAIARTAMA BOTOL CAP	0.14	PAKISTAN	7	0	HEALTH & WELLNESS	
55205	[MILD EXPLOSION] SOAP FOAMING NET BAG MEDIUM	0.14	PERAK	922	20	HEALTH & WELLNESS	
Total rows:	65982						
Total columns:	6						
Category count:							
1. SPORTS & OUTDOORS products							
2. BEAUTY & SKINCARE : 17957 products							
3. HEALTH & WELLNESS : 14015 products							
4. HOME & KITCHEN : 10999 products							
5. WOMEN'S FASHION : 4653 products							
6. MEN'S FASHION : 4653 products							
7. MOTORS & BIKES : 477 products							
Group 3 [Affordable Price]:							

**Figure 17: Output Screenshot of Grouping Product into 4 Categories based on “Price” by using Pandas**

Minimum Number of Total Reviews:	0					
Maximum Number of Total Reviews:	27					
Group 1 (Least popular):						
Product Name	Price	Location	Quantity Sold	Total Reviews	Category	Popularity
12471 BENTON DEEP GREEN TEA CLEANSING FOAM 250g / 10L	28.00	WP KUALA LUMPUR	20	7	BEAUTY & SKINCARE	100.00%
42782 FLAVETTES EFFERVESCENT GLAMZ (105g)	79.00	SELANGOR	29	7	HEALTH & WELLNESS	100.00%
85562 CAPARABA ERASABLE PEN (GOOD-LOOKING)	3.81	N/A	266	7	STATIONERY	100.00%
54996 VINCENCIHAT WALL STICKER MOISTURE-PROOF DECOR	10.40	CHINA	27	7	HOME & LIVING	100.00%
100470 FELTON 10 TIER DRAWER FED 176	76.79	SELANGOR	39	7	STATIONERY	100.00%
2385 BODY MILK WHITENING BRIGHTEN SKIN TONE MOISTUR.	7.89	CHINA	14	7	BEAUTY & SKINCARE	100.00%
92321 (ST/M) 15ML MANG EXTRACT SERUM Pen / PEN. 15ML	25.00	PENANG	22	7	STATIONERY	100.00%
100465 BESSIE WATER MELONADE INHALZ BLACK FRUIT HONEY	5.00	NEGERI SEMBILAN	28	7	STATIONERY	100.00%
12479 VIBRANT GLAMOUR SALICYL ACID ACNE TREATMENT	29.00	CHINA	27	7	BEAUTY & SKINCARE	100.00%
3241 UTX EFFERVESCENT SALTLES SACHET AG 5.28'S	21.90	PAHANG	45	7	BEAUTY & SKINCARE	100.00%
<hr/>						
Total rows:	88729					
Total columns:	7					
<hr/>						
Category count:	10					
- BEAUTY & WELLNESS: 16039 products						
- STATIONERY: 1556 products						
- HOME APPLIANCES: 1553 products						
- PERSONAL CARE: 975 products						
- FOOD & BEVERAGE: 961 products						
- MOTHER & BABY: 2231 products						
<hr/>						
Group 2 (Below Average Popularity):						
Product Name	Price	Location	Quantity Sold	Total Reviews	Category	Popularity
36 REXACOD SAUCER NICOTINAMIDE WHITENING FRECKLE MO.	2.80	WP KUALA LUMPUR	52	14	BEAUTY & SKINCARE	100.00%
111420 EVENING DRESS WOMEN NEW BANGKOK SONGKAIT	198.81	N/A	9	14	WOMEN'S FASHION	100.00%
157 ORS FRACIA CLEANER FACIAL FOAM PINK BIG SIZE	9.00	SELANGOR	44	14	BEAUTY & SKINCARE	100.00%

*Figure 18: Output Screenshot of Grouping Products into 4 Categories based on “Total Reviews” by using Pandas*

	Location	Total_Quantity_Sold	Average_Price	Market Performance
23	SELANGOR	10562173	100.849751	1.065193e+09
14	N/A	18524452	43.647313	8.085426e+08
4	CHINA	6341944	97.622898	6.191190e+08
29	WP KUALA LUMPUR	1227922	101.519881	1.246585e+08
8	JOHOR	1499826	78.075352	1.170567e+08
16	OVERSEAS	2083568	48.480572	1.010126e+08
18	PENANG	962481	100.382696	9.661644e+07
19	PERAK	1029478	80.170026	8.253328e+07
6	HONG KONG	376652	138.300619	5.209120e+07
9	KEDAH	595993	80.050847	4.770974e+07

\*\*\*\*\* Performance \*\*\*\*\*  
Total rows processed: 113596  
Code Execution time: 2.1653 seconds  
Throughput: 52461.36 rows per second  
Current memory usage: 0.8738 MB  
Peak memory usage: 6.2297 MB  
CPU usage: 1% CPU  
\*\*\*\*\*  
CPU times: user 158 ms, sys: 4.42 ms, total: 163 ms  
Wall time: 2.25 s

Figure 19: Output Screenshot of Evaluating and Ranking Market Performance based on “Quantity Sold” for each “Location” by using Pandas

shape: (10, 6)					
Product Name	Price	Location	Quantity Sold	Total Reviews	Category
"Upsee Women's Fashion Heat Res...	17.13	"China"	"753 sold"	"(188)"	"Women's Fashion"
"Summer women's fashionable lig...	13.0	"China"	"54 sold"	"(14)"	"Women's Fashion"
"Summer New Sports Suit Women's...	22.75	"China"	"286 sold"	"(62)"	"Women's Fashion"
"Daring Backless V-Neck Dress S...	19.72	null	"270 sold"	"(37)"	"Women's Fashion"
"Chic Lace Birthday Dress Chri...	12.18	null	"486 sold"	"(33)"	"Women's Fashion"
"LOMOOG Women's Fashion Dress +...	15.23	"China"	"1.6K sold"	"(455)"	"Women's Fashion"
"Summer 2024 Women's Knitted Lo...	22.41	null	"322 sold"	"(26)"	"Women's Fashion"
"Hot Princess Dress Set Sexy Sl...	17.39	null	"92 sold"	"(6)"	"Women's Fashion"
"LOMOOG Women's Fashion Dress +...	14.17	null	"122 sold"	"(42)"	"Women's Fashion"
"Adult Latin Dance Skirt Square...	20.55	null	"350 sold"	"(58)"	"Women's Fashion"

Total rows: 9698  
Total columns: 6

\*\*\*\*\* Performance \*\*\*\*\*  
Total rows processed: 115898  
Code Execution time: 0.6578 seconds  
Throughput: 174031.39 rows per second  
Current memory usage: 2.7796 MB  
Peak memory usage: 2.8168 MB  
CPU usage: 57.8%

Total time for this cell(including time to display the performance):  
CPU times: user 337 ms, sys: 42.8 ms, total: 380 ms  
wall time: 1.07 s

Figure 20: Output Screenshot of Dataset Loading and Display by using Polars

	Product Name	Price	Location	Quantity Sold	Total Reviews	Category
"SAKA GLOWING FOUNDATION SPF ...	3.99	str	164	str	str	str
"Bio-Essence Bio-Gold 24K Radi...	3.5	"Penang"	"55 sold"	"(9)"	"Beauty & Skincare"	
"J'océane Immortelle Divine F...	1.5	"Johor"	"46 sold"	"(10)"	"Beauty & Skincare"	
"LOMOCANE IMMORTELLE DIVINE F...	1.5	"Selangor"	"13 sold"	"(1)"	"Beauty & Skincare"	
"YOBUYI Freckle Cream Effective...	5.45	"China"	"13 sold"	null	"Beauty & Skincare"	
"DT37 Jontam Jolym Nicotimde ...	6.15	"Melaka"	"175 sold"	"(40)"	"Beauty & Skincare"	
"BBHome Jontam Jolym Nicotimde ...	6.19	"Melaka"	"29 sold"	"(0)"	"Beauty & Skincare"	
"DT37 VEZE Men's Roll On 50ML"	5.44	"Melaka"	"129 sold"	"(46)"	"Beauty & Skincare"	
"NEVEA MEN ROLL ON 50ML"	5.0	"WP Kuala Lumpur"	"25 sold"	"(4)"	"Beauty & Skincare"	
"BBHome VEZE Men's Volcanic ...	5.48	"Melaka"	"41 sold"	"(11)"	"Beauty & Skincare"	
"DT37 HYMEYS Facial Cleanser C...	1.53	"Melaka"	"1.3K sold"	"(246)"	"Beauty & Skincare"	

Total rows: 115898  
Total columns: 6

\*\*\*\*\* Performance \*\*\*\*\*  
Total rows processed: 115898  
Code Execution time: 0.1654 seconds  
Throughput: 695962.63 rows per second  
Current memory usage: 8.6284 MB  
Peak memory usage: 9.8989 MB  
CPU usage: 4.08%

Total time for this cell(including time to display the performance):  
CPU times: user 95.5 ms, sys: 47.3 ms, total: 143 ms  
wall time: 1.37 s

Figure 21: Output Screenshot of Dataset Integration by using Polars

shape: (10, 6)					
Product Name	Price	Location	Quantity Sold	Total Reviews	Category
"SAKA GLOWING FOUNDATION SPF ...	3.99	"PENANG"	"55 SOLD"	"(9)"	"BEAUTY & SKINCARE"
"Bio-Essence Bio-Gold 24K RADL..."	3.5	"JOHOR"	"46 SOLD"	"(10)"	"BEAUTY & SKINCARE"
"LOMOCANE IMMORTELLE DIVINE F..."	1.5	"SELANGOR"	"13 SOLD"	"(1)"	"BEAUTY & SKINCARE"
"YOBUYI FRECKLE CREAM EFFECTIVE..."	5.45	"CHINA"	"13 SOLD"	null	"BEAUTY & SKINCARE"
"DT37 JONTAM JOLYUM NICOTIMDE ...	6.15	"MELAKA"	"175 SOLD"	"(40)"	"BEAUTY & SKINCARE"
"BBHOME JONTAM JOLYUM NICOTIMDE ...	6.19	"MELAKA"	"29 SOLD"	"(2)"	"BEAUTY & SKINCARE"
"DT37 VEZE MEN'S VOLCANIC MUD F...	5.44	"MELAKA"	"129 SOLD"	"(46)"	"BEAUTY & SKINCARE"
"NEVEA MEN ROLL ON 50ML"	5.0	"WP KUALA LUMPUR"	"25 SOLD"	"(4)"	"BEAUTY & SKINCARE"
"BBHOME VEZE MEN'S VOLCANIC ..."	5.48	"MELAKA"	"41 SOLD"	"(11)"	"BEAUTY & SKINCARE"
"DT37 HYMEYS FACIAL CLEANSER C...	1.53	"MELAKA"	"1.3K SOLD"	"(246)"	"BEAUTY & SKINCARE"

Total rows: 115898  
Total columns: 6

\*\*\*\*\* Performance \*\*\*\*\*  
Total rows processed: 115898  
Code Execution time: 0.4953 seconds  
Throughput: 158762.39 rows per second  
Current memory usage: 4.0114 MB  
Peak memory usage: 8.0413 MB  
CPU usage: 4.09%

Total time for this cell(including time to display the performance):  
CPU times: user 91.7 ms, sys: 21.8 ms, total: 113 ms  
wall time: 1.08 s

Figure 22: Output Screenshot of Standardization of String Data by using Polars

shape: (10, 6)					
Product Name	Price	Location	Quantity Sold	Total Reviews	Category
str	f64	str	str	int	str
* SAKA GLOWING FOUNDATION SPF ...	3.99	"PENANG"	"55 SOLD"	9	"BEAUTY & SKINCARE"
"BIO-ESSENCE BIO-GOLD 24K RADL...	3.5	"JOHOR"	"48 SOLD"	10	"BEAUTY & SKINCARE"
"LOCCITANE IMMORTELLE DIVINE F...	1.5	"SELANGOR"	"13 SOLD"	1	"BEAUTY & SKINCARE"
"YOUBUY FRECKLE CREAM EFFECTIVE...	5.45	"CHINA"	"13 SOLD"	null	"BEAUTY & SKINCARE"
"D137 JOMTAM JOYVUM NICOTINIDE...	6.15	"MELAKA"	"115 SOLD"	40	"BEAUTY & SKINCARE"
* BHHOME ♥ JOMTAM JOYVUM NICOT...	6.19	"MELAKA"	"29 SOLD"	2	"BEAUTY & SKINCARE"
"D137 VEZE MEN'S VOLCANIC MUD F...	5.44	"MELAKA"	"129 SOLD"	46	"BEAUTY & SKINCARE"
"NEVERA MEN ROLL ON 50ML"	5.0	"WP KUALA LUMPUR"	"25 SOLD"	4	"BEAUTY & SKINCARE"
* BHHOME ♥ VEZE MEN'S VOLCANIC...	5.48	"MELAKA"	"41 SOLD"	11	"BEAUTY & SKINCARE"
"D137 HYMEYS FACIAL CLEANSER C...	1.53	"MELAKA"	"1.3K SOLD"	246	"BEAUTY & SKINCARE"

Total rows: 115699  
Total columns: 6

\*\*\*\*\* Performance \*\*\*\*\*  
Total row processed: 115699  
Code Execution Time: 0.4559 seconds  
Throughput: 248456.72 rows per second  
Current memory usage: 0.0156 MB  
Peak memory usage: 0.0198 MB  
CPU usage: 4.5%

Total time for this cell (including time to display the performance):  
CPU time: 37.7 ms, syst: 1.88 ms, total: 42.8 ms  
wall time: 1.84

Figure 23: Output Screenshot of Converting "Total Reviews" to int Data Type by using Polars

shape: (10, 6)					
Product Name	Price	Location	Quantity Sold	Total Reviews	Category
str	f64	str	str	int	str
* SAKA GLOWING FOUNDATION SPF ...	3.99	"PENANG"	"55"	9	"BEAUTY & SKINCARE"
"BIO-ESSENCE BIO-GOLD 24K RADL...	3.5	"JOHOR"	"46"	10	"BEAUTY & SKINCARE"
"LOCCITANE IMMORTELLE DIVINE F...	1.5	"SELANGOR"	"13"	1	"BEAUTY & SKINCARE"
"YOUBUY FRECKLE CREAM EFFECTIVE...	5.45	"CHINA"	"13"	null	"BEAUTY & SKINCARE"
"D137 JOMTAM JOYVUM NICOTINIDE...	6.15	"MELAKA"	"115"	40	"BEAUTY & SKINCARE"
* BHHOME ♥ JOMTAM JOYVUM NICOT...	6.19	"MELAKA"	"25"	2	"BEAUTY & SKINCARE"
"D137 VEZE MEN'S VOLCANIC MUD F...	5.44	"MELAKA"	"129"	46	"BEAUTY & SKINCARE"
"NEVERA MEN ROLL ON 50ML"	5.0	"WP KUALA LUMPUR"	"25"	4	"BEAUTY & SKINCARE"
* BHHOME ♥ VEZE MEN'S VOLCANIC...	5.48	"MELAKA"	"41"	11	"BEAUTY & SKINCARE"
"D137 HYMEYS FACIAL CLEANSER C...	1.53	"MELAKA"	"1.3K"	246	"BEAUTY & SKINCARE"

Total rows: 115699  
Total columns: 6  
Unique alphabetic characters found: ["K"]

shape: (10, 6)

Product Name	Price	Location	Quantity Sold	Total Reviews	Category
str	f64	str	str	int	str
* SAKA GLOWING FOUNDATION SPF ...	3.99	"PENANG"	55	9	"BEAUTY & SKINCARE"
"BIO-ESSENCE BIO-GOLD 24K RADL...	3.5	"JOHOR"	46	10	"BEAUTY & SKINCARE"
"LOCCITANE IMMORTELLE DIVINE F...	1.5	"SELANGOR"	13	1	"BEAUTY & SKINCARE"
"YOUBUY FRECKLE CREAM EFFECTIVE...	5.45	"CHINA"	13	null	"BEAUTY & SKINCARE"
"D137 JOMTAM JOYVUM NICOTINIDE...	6.15	"MELAKA"	175	40	"BEAUTY & SKINCARE"
* BHHOME ♥ JOMTAM JOYVUM NICOT...	6.19	"MELAKA"	29	2	"BEAUTY & SKINCARE"
"D137 VEZE MEN'S VOLCANIC MUD F...	5.44	"MELAKA"	129	46	"BEAUTY & SKINCARE"
"NEVERA MEN ROLL ON 50ML"	5.0	"WP KUALA LUMPUR"	25	4	"BEAUTY & SKINCARE"
* BHHOME ♥ VEZE MEN'S VOLCANIC...	5.48	"MELAKA"	41	11	"BEAUTY & SKINCARE"
"D137 HYMEYS FACIAL CLEANSER C...	1.53	"MELAKA"	1300	246	"BEAUTY & SKINCARE"

Figure 24: Output Screenshot of Converting "Quantity Sold" to int Data Type by using Polars

Before Handel Missing Values:  
Number of Missing Values for Each Column:

shape: (1, 6)

Product Name_missing	Price_missing	Location_missing	Quantity Sold_missing	Total Reviews_missing	Category_missing
u32	u32	u32	u32	u32	u32
2	0	13517	45566	50958	0

After Handel Missing Values:  
Number of Missing Values for Each Column:

shape: (1, 6)

Product Name_missing	Price_missing	Location_missing	Quantity Sold_missing	Total Reviews_missing	Category_missing
u32	u32	u32	u32	u32	u32
0	0	0	0	0	0

Finalised Dataset:

shape: (10, 6)

Product Name	Price	Location	Quantity Sold	Total Reviews	Category
str	f64	str	str	int	str
* SAKA GLOWING FOUNDATION SPF ...	3.99	"PENANG"	"55"	9	"BEAUTY & SKINCARE"
"BIO-ESSENCE BIO-GOLD 24K RADL...	3.5	"JOHOR"	"46"	10	"BEAUTY & SKINCARE"
"LOCCITANE IMMORTELLE DIVINE F...	1.5	"SELANGOR"	"13"	1	"BEAUTY & SKINCARE"

Figure 25: Output Screenshot of Checking and Handling Missing Values by using Polars

Product Name	Price	Location	Quantity Sold	Total Reviews	Category
str	f64	str	str	int	str
"YONTOO ) HADA LABO GOKUEN PR...	50.19	"JAPAN"	0	1	"BEAUTY & SKINCARE"
"YONTOO ) HADA LABO GOKUEN PR...	50.19	"JAPAN"	0	1	"BEAUTY & SKINCARE"
"11 PCS 10X7 3D WALLPAPER BLI...	1.98	"PERAK"	3400	93	"HOME & LIVING"
* READY STOCK ❸ 超便宜壁紙促销品项! ...	47.6	"WP KUALA LUMPUR"	8	5	"BEAUTY & SKINCARE"
* STOCK SODA ADAB # TAMANERA BLA...	22.0	"PENANG"	0	0	"BEAUTY & SKINCARE"
* STOCK SODA ADAB # TAMANERA BLA...	22.0	"PENANG"	0	0	"BEAUTY & SKINCARE"
* COCONUT OIL NATURAL LIP BALM	12.0	"SELANGOR"	0	0	"BEAUTY & SKINCARE"
* COCONUT OIL NATURAL LIP BALM	12.0	"SELANGOR"	0	0	"BEAUTY & SKINCARE"

Total rows: 2854  
Total columns: 6

After Handling Duplicate:  
Number of Duplicate Rows: 0

Product Name	Price	Location	Quantity Sold	Total Reviews	Category
str	f64	str	str	int	str
"NEW PROMO FRESHFRESH COMFORT EX...	45.86	"JOHOR"	0	0	"BEAUTY & SKINCARE"
"NEW YEAR CARDBOARD CAT LUCKY CAT...	21.7	"KUL"	83	10	"HOME & LIVING"
"EEN CAPSULE WHITENING / EEN BQ...	39.0	"PERAK"	6	3	"HEALTH & WELLNESS"

Total rows: 8  
Total columns: 6

Finalised Dataset:

shape: (10, 6)

Performance					
Total rows processed: 113596					
Code Execution time: 0.0959 seconds					
Throughput: 1184103.69 rows per second					
Current memory usage: 0.0081 MB					
Peak memory usage: 0.0143 MB					
CPU usage: 4.0%					
=====					
Total time for this cell(Including time to display the performance):					
CPU times: user 112 ms, sys: 33.7 ms, total: 145 ms					
Wall time: 1.1 s					

Figure 27: Output Screenshot of Exporting Cleaned Dataset File by using Polars

Min Price: 0.05 Max Price: 172.82					
Group 1 (Budget Friendly Price): 65932 products					
Group 2 (Affordable Price): 22826 products					
Group 3 (Mid-Range Price): 18027 products					
Group 4 (Premium Price): 14603 products					
Group 1 (Budget Friendly Price):					
shape: (10, 6)					
Product Name	Price	Location	Quantity Sold	Total Reviews	Category
"MALAYSIA 3PIN PLUG"	0.05	"SELANGOR"	0	0	"HOME APPLIANCES"
"LIVE TRACKING PINOT FOR SALES"	0.05	"KELANTAN"	0	1	"HOME & LIVING"
"SPOT BESEN PLUS EYE CARE + BR..."	0.1	"JOHOR"	0	0	"HEALTH & WELLNESS"
"1PCS 0.5MM BALL GEL INK PEN MA..."	0.1	"SELANGOR"	33500	135	"STATIONERY"
"HEALTH TREE RESCUE PRODUCT U..."	0.1	"CHINA"	10	3	"HEALTH & WELLNESS"
"BEFREE BESEN PLUS VITAMIN EYE..."	0.1	"JOHOR"	0	0	"HEALTH & WELLNESS"
"READY STOCK MICKEY MOUSE MASCO..."	0.11	"WP KUALA LUMPUR"	49	1	"HOME & LIVING"
"HARI RAYA HUSAN GANTUNG KRA..."	0.12	"WP KUALA LUMPUR"	24	1	"HOME & LIVING"
"PINUTUP BOTTLE CAPS ARTIMARAK..."	0.13	"PENANG"	7	0	"HEALTH & WELLNESS"
"(MILD EXFOLIATION) SOAP FOAM..."	0.14	"PERAK"	92	20	"HEALTH & WELLNESS"
Total rows: 65932					
Total columns: 6					
■ Category count:					
- STATIONERY: 22826 products					
- HEALTH & WELLNESS: 18027 products					
- HOME & LIVING: 22936 products					
- MOTHER & BABY: 14603 products					
- HOME APPLIANCES: 3887 products					
- MOTHER & BABY: 477 products					
Group 2 (Affordable Price):					

Figure 28: Output Screenshot of Grouping Product into 4 Categories based on "Price" by using Polars

Minimum Number of Total Reviews: 0 Maximum Number of Total Reviews: 27					
Group 1 (Least popular): 88729 products					
Group 2 (Below Average Popularity): 7846 products					
Group 3 (Above Average Popularity): 4653 products					
Group 4 (Most popular): 2783 products					
Group 1 (Least popular):					
Product Name	Price	Location	Quantity Sold	Total Reviews	Category
"UNICORN STATIONERY 0.5 MM RUB..."	2.0	"SELANGOR"	68	7	"STATIONERY"
"HOME LIVING ROOM BEDROOM FLOOR..."	55.0	"NEGERI SEMBILAN"	28	7	"HOME & LIVING"
"ELBA SOL ELECTRIC KETTLE STA..."	119.0	"WP KUALA LUMPUR"	19	7	"HOME APPLIANCES"
"BURTS BEES 100% NATURAL MOSTU..."	35.9	"PENANG"	44	7	"BEAUTY & SKINCARE"
"SOLID ADHESIVE NAIL GLUE SUPER..."	2.8	"MELAKA"	31	7	"HEALTH & WELLNESS"
"SHIP 24H 8 PCS CURVING WRITING..."	10.3	"CHINA"	16	7	"STATIONERY"
"QUICK EXTENDOID GLUE MANICURE A..."	4.56	"N/A"	81	7	"STATIONERY"
"(笔)(3mm) 罗经笔; GOLD PEN / 1支..."	3.5	"SELANGOR"	37	7	"STATIONERY"
"GLUMONY KAPSUL PEMUTIH BADAN W..."	24.03	"WP KUALA LUMPUR"	20	7	"HEALTH & WELLNESS"
Total rows: 113596					
Total columns: 6					
=====					
Total time for this cell(Including time to display the performance):					
CPU times: user 42 ms, sys: 15.1 ms, total: 57.1 ms					
Wall time: 1.11 s					

Figure 29: Output Screenshot of Grouping Products into 4 Categories based on "Total Reviews" by using Polars

Performance					
Total rows processed: 113596					
Code Execution time: 0.1111 seconds					
Throughput: 101119.21 rows per second					
Current memory usage: 0.0098 MB					
Peak memory usage: 0.0222 MB					
CPU usage: 100.0%					
=====					
Total time for this cell(Including time to display the performance):					
CPU times: user 42 ms, sys: 15.1 ms, total: 57.1 ms					
Wall time: 1.11 s					

Figure 30: Output Screenshot of Evaluating and Ranking Market Performance based on "Quantity Sold" for each "Location" by using Polars

Product Name	Price	Location	Quantity Sold	Total Reviews	Category
0 Upsee Women's Fashion Heat Resistant Long Crl...	17.13	China	753 sold	(1B8)	Women's Fashion
1 Summer women's fashionable high-end loose and ...	13.00	China	54 sold	(14)	Women's Fashion
2 Summer New Sports Suit Women's Fashion Slim Pr...	22.75	China	286 sold	(62)	Women's Fashion
3 Daring Backless V-Neck Dress Short Mini Skirt ...	19.72	None	270 sold	(37)	Women's Fashion
4 Chic Lace Birthday Dress Christmas Loose Sessa...	12.18	None	486 sold	(35)	Women's Fashion
5 LOMOGI Women's Fashion Dress + Outer Cardigan ...	15.23	China	1.6K sold	(455)	Women's Fashion
6 Summer 2024 Women's Knitted Long Hollow Out Ve...	22.41	None	322 sold	(26)	Women's Fashion
7 Hot Princess Dress Set Sexy Slimming Waist Bod...	17.39	None	92 sold	(6)	Women's Fashion
8 LOMOGI Women's Fashion Dress + Outer CardigSite...	14.17	None	122 sold	(42)	Women's Fashion
9 Adult Latin Dance Skirt Square Dance Costume H...	20.55	None	350 sold	(58)	Women's Fashion

Total rows: 9698  
Total columns: 6

===== Performance =====

Total rows processed: 113998  
Code Execution time: 133.0075 seconds  
Throughput: 8018.43 rows per second  
Current memory usage: 10.9488 MB  
Peak memory usage: 34.7847 MB  
CPU usage: 84.8%

Total time for this cell (Including time to display the performance):  
CPU times: user 1min 31s, sys: 469 ms, total: 1min 32s  
Wall time: 1min 54s

Product Name	Price	Location	Quantity Sold	Total Reviews	Category
0 ❤ SAKA GLOWING FOUNDATION SPF 50+ / FW Fatin W...	3.99	PENANG	55 SOLD	(9)	BEAUTY & SKINCARE
1 BIO-ESSENCE BIO-GOLD 24K RADIANCE/WHITEISING/B...	3.50	JOHOR	46 SOLD	(10)	BEAUTY & SKINCARE
2 L'Occitane Immortelle Divine Foaming Cleansing...	1.50	SELANGOR	13 SOLD	(1)	BEAUTY & SKINCARE
3 YOUNBUY FRECKLE CREAM EFFECTIVELY REMOVE MELASMA...	9.45	CHINA	13 SOLD	NAN	BEAUTY & SKINCARE
4 DT37 JONITAM JOLYUM NICOTINEME AMINO ACID FACIA...	6.15	MELAKA	175 SOLD	(40)	BEAUTY & SKINCARE
5 ❤BBHome ❤ JONITAM JOLYUM NICOTINEME AMINO ACID...	6.19	MELAKA	29 SOLD	(2)	BEAUTY & SKINCARE
6 DT37 VEZE MEN'S VOLCANIC MUD FACIAL CLEANSER B...	5.44	MELAKA	129 SOLD	(46)	BEAUTY & SKINCARE
7 NEVEA MEN ROLL ON 500ML 5.00 WP KUALA LUMPUR	5.00	WP KUALA LUMPUR	25 SOLD	(4)	BEAUTY & SKINCARE
8 ❤BBHome ❤ VEZE MEN'S VOLCANIC MUD FACIAL CLEA...	5.48	MELAKA	41 SOLD	(11)	BEAUTY & SKINCARE
9 DT37 HYMEY'S FACIAL CLEANSER CREAM WHITENING M...	1.53	MELAKA	1.3K SOLD	(246)	BEAUTY & SKINCARE

Total rows: 115988  
Total columns: 6

===== Performance =====

Total rows processed: 115,098  
Code Execution time: 4.0765 seconds  
Throughput: 21,637 rows per second  
Current memory usage: 6.0738 MB  
Peak memory usage: 6.1984 MB  
CPU usage: 12.6%

Total time for this cell (Including time to display the performance):  
CPU times: user 176 ms, sys: 7 ms, total: 183 ms  
Wall time: 5.88 s

Figure 31: Output Screenshot of Dataset Loading and Display by using PySpark

Product Name	Price	Location	Quantity Sold	Total Reviews	Category
0 ❤ SAKA GLOWING FOUNDATION SPF 50+ / FW Fatin W...	3.99	PENANG	55 SOLD	(9)	BEAUTY & SKINCARE
1 BIO-ESSENCE BIO-GOLD 24K RADIANCE/WHITEISING/B...	3.50	JOHOR	46 SOLD	(10)	BEAUTY & SKINCARE
2 L'Occitane Immortelle Divine Foaming Cleansing...	1.50	SELANGOR	13 SOLD	(1)	BEAUTY & SKINCARE
3 YOUNBUY FRECKLE CREAM EFFECTIVELY REMOVE MELASMA...	9.45	CHINA	13 SOLD	NAN	BEAUTY & SKINCARE
4 DT37 JONITAM JOLYUM NICOTINEME AMINO ACID FACIA...	6.15	MELAKA	175 SOLD	(40)	BEAUTY & SKINCARE
5 ❤BBHome ❤ JONITAM JOLYUM NICOTINEME AMINO ACID...	6.19	MELAKA	29 SOLD	(2)	BEAUTY & SKINCARE
6 DT37 VEZE MEN'S VOLCANIC MUD FACIAL CLEANSER B...	5.44	MELAKA	129 SOLD	(46)	BEAUTY & SKINCARE
7 NEVEA MEN ROLL ON 500ML 5.00 WP KUALA LUMPUR	5.00	WP KUALA LUMPUR	25 SOLD	(4)	BEAUTY & SKINCARE
8 ❤BBHome ❤ VEZE MEN'S VOLCANIC MUD FACIAL CLEA...	5.48	MELAKA	41 SOLD	(11)	BEAUTY & SKINCARE
9 DT37 HYMEY'S FACIAL CLEANSER CREAM WHITENING M...	1.53	MELAKA	1.3K SOLD	(246)	BEAUTY & SKINCARE

Total rows: 115990  
Total columns: 6

===== Performance =====

Total rows processed: 115,090  
Code Execution time: 4.2237 seconds  
Throughput: 22,356.57 rows per second  
Current memory usage: 8.0679 MB  
Peak memory usage: 8.1824 MB  
CPU usage: 61.0%

Total time for this cell (Including time to display the performance):  
CPU times: user 205 ms, sys: 12 ms, total: 217 ms  
Wall time: 5.21 s

Figure 33: Output Screenshot of Standardization of String Data by using PySpark

Product Name	Price	Location	Quantity Sold	Total Reviews	Category
0 ❤ SAKA GLOWING FOUNDATION SPF 50+ / FW Fatin W...	3.99	PENANG	55	9.0	BEAUTY & SKINCARE
1 BIO-ESSENCE BIO-GOLD 24K RADIANCE/WHITEISING/B...	3.50	JOHOR	46	10.0	BEAUTY & SKINCARE
2 L'Occitane Immortelle Divine Foaming Cleansing...	1.50	SELANGOR	13	1.0	BEAUTY & SKINCARE
3 YOUNBUY FRECKLE CREAM EFFECTIVELY REMOVE MELASMA...	9.45	CHINA	13	NAN	BEAUTY & SKINCARE
4 DT37 JONITAM JOLYUM NICOTINEME AMINO ACID FACIA...	6.15	MELAKA	175	40.0	BEAUTY & SKINCARE
5 ❤BBHome ❤ JONITAM JOLYUM NICOTINEME AMINO ACID...	6.19	MELAKA	29	2.0	BEAUTY & SKINCARE
6 DT37 VEZE MEN'S VOLCANIC MUD FACIAL CLEANSER B...	5.44	MELAKA	129	46.0	BEAUTY & SKINCARE
7 NEVEA MEN ROLL ON 500ML 5.00 WP KUALA LUMPUR	5.00	WP KUALA LUMPUR	25	4.0	BEAUTY & SKINCARE
8 ❤BBHome ❤ VEZE MEN'S VOLCANIC MUD FACIAL CLEA...	5.48	MELAKA	41	11.0	BEAUTY & SKINCARE
9 DT37 HYMEY'S FACIAL CLEANSER CREAM WHITENING M...	1.53	MELAKA	1300	246.0	BEAUTY & SKINCARE

Total rows: 115990  
Total columns: 6

===== Performance =====

Total rows processed: 115,090  
Code Execution time: 4.1328 seconds  
Throughput: 22,356.88 rows per second  
Current memory usage: 8.4972 MB  
Peak memory usage: 8.1938 MB  
CPU usage: 6.4%

Total time for this cell (Including time to display the performance):  
CPU times: user 213 ms, sys: 0.99 ms, total: 222 ms  
Wall time: 5.14 s

Figure 35: Output Screenshot of Converting “Quantity Sold” to int Data Type by using PySpark

Product Name	Price	Location	Quantity Sold	Total Reviews	Category
0 ❤ SAKA GLOWING FOUNDATION SPF 50+ / FW Fatin W...	3.99	PENANG	55 sold	(9)	Beauty & Skincare
1 Bio-Essence Bio-Gold 24k Radiance/Whitening/B...	3.50	JOHOR	46 sold	(10)	Beauty & Skincare
2 L'occitane Immortelle Divine Foaming Cleansing...	1.50	SELANGOR	13 sold	(1)	Beauty & Skincare
3 YOUNBUY Freckle Cream Effectively Remove Melasma...	5.45	CHINA	13 sold	Nan	Beauty & Skincare
4 DT37 Jonitam Jolyum Nicotinamide Amino Acid Facia...	6.15	MELAKA	175 sold	(40)	Beauty & Skincare
5 ❤BBHome ❤ Jonitam Jolyum Nicotinamide Amino Acid...	6.19	MELAKA	29 sold	(2)	Beauty & Skincare
6 DT37 VEZE Men's Volcanic Mud Facial Cleanser B...	5.44	MELAKA	129 sold	(46)	Beauty & Skincare
7 NEVEA Men Roll On 500ML 5.00 WP KUALA LUMPUR	5.00	WP KUALA LUMPUR	25 sold	(4)	Beauty & Skincare
8 ❤BBHome ❤ VEZE Men's Volcanic Mud Facial Clea...	5.48	MELAKA	41 sold	(11)	Beauty & Skincare
9 DT37 HYMEY'S Facial Cleanser Cream Whitening M...	1.53	MELAKA	1.3K sold	(246)	Beauty & Skincare

Total rows: 115998  
Total columns: 6

===== Performance =====

Total rows processed: 115,098  
Code Execution time: 4.0765 seconds  
Throughput: 21,637 rows per second  
Current memory usage: 6.0738 MB  
Peak memory usage: 6.1984 MB  
CPU usage: 12.6%

Total time for this cell (Including time to display the performance):  
CPU times: user 176 ms, sys: 7 ms, total: 183 ms  
Wall time: 5.88 s

Figure 32: Output Screenshot of Dataset Integration by using PySpark

Product Name	Price	Location	Quantity Sold	Total Reviews	Category
0 ❤ SAKA GLOWING FOUNDATION SPF 50+ / FW Fatin W...	3.99	PENANG	55 SOLD	9.0	BEAUTY & SKINCARE
1 BIO-ESSENCE BIO-GOLD 24K RADIANCE/WHITEISING/B...	3.50	JOHOR	46 SOLD	10.0	BEAUTY & SKINCARE
2 L'OCCTANE IMMORTELLE DIVINE FOAMING CLEANSING...	1.50	SELANGOR	13 SOLD	1.0	BEAUTY & SKINCARE
3 YOUNBUY FRECKLE CREAM EFFECTIVELY REMOVE MELASMA...	5.45	CHINA	13 SOLD	NAN	BEAUTY & SKINCARE
4 DT37 JONITAM JOLYUM NICOTINEME AMINO ACID FACIA...	6.15	MELAKA	175 SOLD	40.0	BEAUTY & SKINCARE
5 ❤BBHome ❤ JONITAM JOLYUM NICOTINEME AMINO ACID...	6.19	MELAKA	29 SOLD	2.0	BEAUTY & SKINCARE
6 DT37 VEZE MEN'S VOLCANIC MUD FACIAL CLEANSER B...	5.44	MELAKA	129 SOLD	46.0	BEAUTY & SKINCARE
7 NEVA MEN ROLL ON 500ML 5.00 WP KUALA LUMPUR	5.00	WP KUALA LUMPUR	25 SOLD	4.0	BEAUTY & SKINCARE
8 ❤BBHome ❤ VEZE MEN'S VOLCANIC MUD FACIAL CLEA...	5.48	MELAKA	41 SOLD	11.0	BEAUTY & SKINCARE
9 DT37 HYMEY'S FACIAL CLEANSER CREAM WHITENING M...	1.53	MELAKA	1.3K SOLD	246.0	BEAUTY & SKINCARE

Total rows: 115998  
Total columns: 6

===== Performance =====

Total rows processed: 115,098  
Code Execution time: 4.0765 seconds  
Throughput: 21,637 rows per second  
Current memory usage: 6.0738 MB  
Peak memory usage: 6.1984 MB  
CPU usage: 12.6%

Total time for this cell (Including time to display the performance):  
CPU times: user 176 ms, sys: 7 ms, total: 183 ms  
Wall time: 5.88 s

Figure 34: Output Screenshot of Converting “Total Reviews” to int Data Type by using PySpark

Product Name	Price	Location	Quantity Sold	Total Reviews	Category
6 DT37 VEZE MEN'S VOLCANIC MUD FACIAL CLEANSER B...	5.44	MELAKA	129	46.0	BEAUTY & SKINCARE
7 NEVEA MEN ROLL ON 500ML 5.00 WP KUALA LUMPUR	5.00	WP KUALA LUMPUR	25	4.0	BEAUTY & SKINCARE
8 ❤BBHome ❤ VEZE MEN'S VOLCANIC MUD FACIAL CLEA...	5.48	MELAKA	41	11.0	BEAUTY & SKINCARE
9 DT37 HYMEY'S FACIAL CLEANSER CREAM WHITENING M...	1.53	MELAKA	1300	246.0	BEAUTY & SKINCARE

Before Handling Missing Values

Missing Values:

Product Name_missing	Price_missing	Location_missing	Quantity Sold_missing	Total Reviews_missing	Category_missing
0	2	0	13517	45566	0

After Handling Missing Values

Missing Values:

Product Name_missing	Price_missing	Location_missing	Quantity Sold_missing	Total Reviews_missing	Category_missing
0	0	0	0	0	0

Finalised Dataset:

Product Name	Price	Location	Quantity Sold	Total Reviews	Category
0 ❤ SAKA GLOWING FOUNDATION SPF 50+ / FW Fatin W...	3.99	PENANG	55	9	BEAUTY & SKINCARE
1 BIO-ESSENCE BIO-GOLD 24K RADIANCE/WHITEISING/B...	3.50	JOHOR	46	10	BEAUTY & SKINCARE
2 L'OCCTANE IMMORTELLE DIVINE FOAMING CLEANSING...	1.50	SELANGOR	13	1	BEAUTY & SKINCARE
3 YOUNBUY FRECKLE CREAM EFFECTIVELY REMOVE MELASMA...	5.45	CHINA	13	0	BEAUTY & SKINCARE
4 DT37 JONITAM JOLYUM NICOTINEME AMINO ACID FACIA...	6.15	MELAKA	175	40	BEAUTY & SKINCARE
5 ❤BBHome ❤ JONITAM JOLYUM NICOTINEME AMINO ACID...	6.19	MELAKA	29	2	BEAUTY & SKINCARE
6 DT37 VEZE MEN'S VOLCANIC MUD FACIAL CLEANSER B...	5.44	MELAKA	129	46	BEAUTY & SKINCARE
7 NEVEA MEN ROLL ON 500ML 5.00 WP KUALA LUMPUR	5.00	WP KUALA LUMPUR	25	4	BEAUTY & SKINCARE

**Figure 36: Output Screenshot of Checking and Handling Missing Values by using PySpark**

3	( ROHTO ) HADA-LABO GOKUIN PREMIUM HYALURONIC... 50.19	JAPAN	0	1	BEAUTY & SKINCARE
4	(1 PCS) 70K77 3D WALLPAPER BRICK WALL STICKERS... 1.98	PERAK	3400	93	HOME & LIVING
-	-	-	-	-	-
2849	READY STOCK 6 PROMOTION 100% ORIGINAL 美乐家...	47.00	WP KUALA LUMPUR	8	5 BEAUTY & SKINCARE
2850	STOCK SEDNA ADA TANAMERA BLACK FORMULATION PA...	22.00	PENANG	0	0 BEAUTY & SKINCARE
2851	STOCK SEDNA ADA TANAMERA BLACK FORMULATION PA...	22.00	PENANG	0	0 BEAUTY & SKINCARE
2852	COCONUT OIL NATURAL LIP BALM ( Revived... 12.00	SELANGOR	0	0 BEAUTY & SKINCARE	
2853	COCONUT OIL NATURAL LIP BALM ( Revived... 12.00	SELANGOR	0	0 BEAUTY & SKINCARE	

2854 rows × 6 columns

Total rows: 2854

Total columns: 6

After Handling Duplicates

Number of Duplicate rows: 0

Product Name Price Location Quantity Sold Total Reviews Category

Total rows: 0

Total columns: 6

Finalised Dataset:

	Product Name	Price	Location	Quantity Sold	Total Reviews	Category
0	ORIGINAL GLOW FOUNDATION BY HQ EKIN BEAUTY SPF60	8.29	PERAK	0	1 BEAUTY & SKINCARE	
1	BEAUTY FORMULAS MAKE UP REMOVER CLEANSING FAC...	8.90	JOHOR	592	93 BEAUTY & SKINCARE	
2	BIGAQUA PAPAYA CLEANSING WITH VITAMINS 100 GRAM	10.00	WP KUALA LUMPUR	49	6 BEAUTY & SKINCARE	
3	FACIAL CLEANSER WHITENING AND FRECKLE REMOVING...	7.39	SELANGOR	484	132 BEAUTY & SKINCARE	
4	ALOE VERA 99% SOOTHING GEL LIPSTICK LIP BALM	8.50	WP KUALA LUMPUR	178	44 BEAUTY & SKINCARE	
5	20G CHEILITIS CREAM LIP CARE CHEILITIS REPAIR ...	9.14	CHINA	31	8 BEAUTY & SKINCARE	
6	SNEFE WHITE LILY HYDRATING CLEANSER 雪特絲滋潤清潔乳	12.00	PENANG	47	20 BEAUTY & SKINCARE	
7	LIP PLUMP SERUM INCREASE LIP ELASTICITY REDUCE...	7.61	CHINA	6	1 BEAUTY & SKINCARE	
8	SABUN GLOW GLOWING READY STOCK	9.20	KELANTAN	0	0 BEAUTY & SKINCARE	
9	SAFI YOUTH GOLD SERIES (FACIAL CLEANSER / EXFO... 10.41	PERAK	0	0 BEAUTY & SKINCARE		

**Figure 37: Output Screenshot of Checking and Handling Duplicates by using PySpark**

```
=====
Performance
=====
Total rows processed: 113596
Code Execution time: 10.6576 seconds
Throughput: 18658.73 rows per second
Current memory usage: 0.0334 MB
Peak memory usage: 0.0888 MB
CPU usage: 75.1%
```

Total time for this cell(Including time to display the performance):  
CPU times: user 148 ms, sys: 8.15 ms, total: 148 ms  
Wall time: 11.7 s

**Figure 38: Output Screenshot of Exporting Cleaned Dataset File by using PySpark**

Group 1 (Budget Friendly Price):
Product Name Price Location Quantity Sold Total Reviews Category
0 ORIGINAL GLOW FOUNDATION BY HQ EKIN BEAUTY SPF60 8.29 PERAK 0 1 BEAUTY & SKINCARE
1 BEAUTY FORMULAS MAKE UP REMOVER CLEANSING FAC... 8.90 JOHOR 592 93 BEAUTY & SKINCARE
2 BIGAQUA PAPAYA CLEANSING WITH VITAMINS 100 GRAM 10.00 WP KUALA LUMPUR 49 6 BEAUTY & SKINCARE
3 FACIAL CLEANSER WHITENING AND FRECKLE REMOVING... 7.39 SELANGOR 484 132 BEAUTY & SKINCARE
4 ALOE VERA 99% SOOTHING GEL LIPSTICK LIP BALM 8.50 WP KUALA LUMPUR 178 44 BEAUTY & SKINCARE
5 20G CHEILITIS CREAM LIP CARE CHEILITIS REPAIR ... 9.14 CHINA 31 8 BEAUTY & SKINCARE
6 SNEFE WHITE LILY HYDRATING CLEANSER 雪特絲滋潤清潔乳 12.00 PENANG 47 20 BEAUTY & SKINCARE
7 LIP PLUMP SERUM INCREASE LIP ELASTICITY REDUCE... 7.61 CHINA 6 1 BEAUTY & SKINCARE
8 SABUN GLOW GLOWING READY STOCK
9 SAFI YOUTH GOLD SERIES (FACIAL CLEANSER / EXFO... 10.41 PERAK 0 0 BEAUTY & SKINCARE

Total rows: 65992

Total columns: 6

- Category count:
  - STATOERY: 20468 products
  - BEAUTY & SKINCARE: 17975 products
  - HEALTH & MEDICAL: 1035 products
  - HOME & LIVING: 8293 products
  - HOME APPLIANCES: 3837 products
  - OTHER & BABY: 477 products
- Group 2 (Affordable Price):

**Figure 39: Output Screenshot of Grouping Product into 4 Categories based on “Price” by using PySpark**

Filtered minimum: 0
Filtered maximum: 27
Group 1 (Least popular): 88729 products
Group 2 (Below Average Popularity): 7945 products
Group 3 (Above Average Popularity): 4653 products
Group 4 (Most Popular): 2787 products
Group 1 (Least popular):
Product Name Price Location Quantity Sold Total Reviews Category
0 ORIGINAL GLOW FOUNDATION BY HQ EKIN BEAUTY SPF60 8.29 PERAK 0 1 BEAUTY & SKINCARE
1 BIGAQUA PAPAYA CLEANSING WITH VITAMINS 100 GRAM 10.00 WP KUALA LUMPUR 49 6 BEAUTY & SKINCARE
2 LIP PLUMP SERUM INCREASE LIP ELASTICITY REDUCE... 7.61 CHINA 6 1 BEAUTY & SKINCARE
3 SABUN GLOW GLOWING READY STOCK 9.20 KELANTAN 0 0 BEAUTY & SKINCARE
4 SAFI YOUTH GOLD SERIES (FACIAL CLEANSER / EXFO... 10.41 PERAK 0 0 BEAUTY & SKINCARE
5 [HEALTHCARE.ONLINE PHARMACY] IF SULFUR SKIN SO... 8.80 JOHOR 0 0 BEAUTY & SKINCARE
6 [CLEARANCE] HIMALAYA MOISTURISING ALOE VERA FA... 8.30 SELANGOR 33 7 BEAUTY & SKINCARE
7 JOJI SRN BUBBLE SOAP 9.90 PENANG 0 0 BEAUTY & SKINCARE
8 PACKAGING BARU ! WILYA WHITENING SOAP / SABUN ... 10.90 KELANTAN 0 0 BEAUTY & SKINCARE
9 FLACENTA UV-WHITENING HAND & BODY LOTION 9.90 SELANGOR 14 5 BEAUTY & SKINCARE

**Figure 40: Output Screenshot of Grouping Products into 4 Categories based on “Total Reviews” by using PySpark**

	Location	Total Quantity Sold	Average Price	Market Performance
0	SELANGOR	10562173	100.849751	1.065193e+09
1	N/A	18524452	43.647313	8.085426e+08
2	CHINA	6341944	97.622898	6.191190e+08
3	WP KUALA LUMPUR	1227922	101.519881	1.246585e+08
4	JOHOR	1499826	78.073532	1.170967e+08
5	OVERSEAS	2083568	48.480572	1.010126e+08
6	PENANG	962481	100.382696	9.661644e+07
7	PERAK	1029478	80.170026	8.253328e+07
8	HONG KONG	376652	138.300619	5.209120e+07
9	KEDAH	595993	80.050847	4.770974e+07

===== Performance =====

```
Total rows processed: 113,598
Code Execution time: 1.8099 seconds
Throughput: 62,763.71 rows per second
Current memory usage: 0.0618 MB
Peak memory usage: 0.1632 MB
CPU usage: 56.9% 

Total time for this cell (Including time to display the performance):
CPU times: user 148 ms, sys: 5.9 ms, total: 154 ms
Wall time: 2.82 s
```

**Figure 41: Output Screenshot of Evaluating and Ranking Market Performance based on “Quantity Sold” for each “Location” by using PySpark**

## Links to full code repo or dataset

### 1. Pandas:

- [Part 1](#)
- [Part 2](#)

### 2. Polars:

- [Part 1](#)
- [Part 2](#)

### 3. PySpark:

- [Part 1](#)
- [Part 2](#)

### 4. Dataset

- [Raw Dataset](#)
- [Cleaned Dataset](#)

# HPDP\_Project Report.pdf

## ORIGINALITY REPORT



## PRIMARY SOURCES

---

1	<b>neptune.ai</b> Internet Source	1 %
2	<b>www.diva-portal.org</b> Internet Source	<1 %
3	<b>2captcha.com</b> Internet Source	<1 %
4	<b>Submitted to University of Surrey</b> Student Paper	<1 %
5	<b>Submitted to Universiti Teknologi Malaysia</b> Student Paper	<1 %
6	<b>documents1.worldbank.org</b> Internet Source	<1 %
7	<b>Submitted to Arab Open University</b> Student Paper	<1 %
8	<b>Arvind Dagur, Karan Singh, Pawan Singh Mehra, Dhirendra Kumar Shukla. "Intelligent Computing and Communication Techniques - Volume 2", CRC Press, 2025</b> Publication	<1 %

---

---

Exclude quotes      On

Exclude bibliography      On

Exclude matches      Off