# Analysing Online Retail Transactions using Big Data Framework

1 author:

Aashish Prasad
National College of Ireland
**8** PUBLICATIONS   **6** CITATIONS

# Analysing Online Retail Transactions using Big Data Framework

Aashish Prasad

*X17170826*

*MSc in Data Analytics*

*National College of Ireland*

Dublin, Ireland

*28th April 2019*

*Abstract*—Easy access to the internet has broadened the use of online retail and shopping websites. An analysis of the enormous amount of data generated from such transactions is important as part of business growth and development. In this research, we used a distributed computing approach to analyse big data using Hadoop. The dataset was processed and loaded into MySQL as an input for other data processing frameworks. We used sqoop for data transfer from MySQL into HDFS. As per the business objectives, the tasks were distributed among Java, Pig, Hive and Spark frameworks for which the results were generated and stored into HBase database. Further, the outputs were visualized to provide graphical business insights.

*Keywords*— Hadoop, MapReduce, Sqoop, Pig, Hive, Spark, HBase, MySQL, HDFS

## I. INTRODUCTION

In the era of internet, data growth is a constant process. This has opened better opportunities for consumer market analysis in terms of quality with larger samples. Many organizations involved in direct and indirect interaction with the consumer market have realized the potential of data analysis of such large datasets for business growth and expansion. With the growing number of internet access to people, the consumer demand for products and services have also increased over the online retail and shopping websites.[1]

To meet such business requirements, many approaches have been used proposed using traditional data management techniques, such as relational database management system, which often result in low system performance due to architectural incompetence. Here, comes Hadoop, a distributed computing framework that follows MapReduce programming model for handling large dataset and well known for its fault tolerance feature [2]. It provides the advantage of high-speed processing with its data splitting technique.

This research aims to analyse online retail transaction dataset using MapReduce implementation of Java, Pig, Hive and Spark. Also, technologies such as Hadoop Distributed File System (HDFS), Sqoop, MySQL and HBase has been used for data storage and processing. The business objectives and research questions are elaborated in the next section.

### A. Objective and Research Questions

The objective of the research is to find meaningful insights that can be useful in the development of an online retail store. This objective is listed in the following Business queries as mentioned below.

1) What are the top 20 most popular products among customers from the United Kingdom, sold in the first quarter of the year?
2) What is the average number of online orders received on an hourly basis in a week?
3) How much was Net Profit on a Quarterly basis for the year 2011?
4) What is the average order amount for each country?
5) In which month the sales revenue reached a maximum?
6) What is the total number of customers from each country?
7) Which top 10 customer ID purchased the most?

### B. Motivation for Research

According to [3], the online retail market is price sensitive. Promotional offers and discounts offered to the customers have a high impact on the revenue in this sector. Also, it is necessary to understand the consumer demands for products to maintain stock availability or avoid the unnecessary stockpile up. Finding the correct answer to such critical business queries is only possible through data analytic approach. Hence, a system to analyse big data is an essential step towards business goals. This system will be useful in making correct business decision and strategies.

Section II discusses, some of the previous work and research accomplished using similar technologies that we have used in this research. While Section III, details of dataset and methodology used for prepossessing and overall system development process has been explained. Results for the study using data visualisation techniques has been presented in Section IV.

## II. RELATED WORK

To analyse large data on internet traffic, research was conducted by [4], using MapReduce based frameworks. According to the author, MapReduce processes

data efficiently with its distributed parallel processing mechanism that includes a map and reduce functions. To analyse ever growing internet traffic data, a traffic analysis system using Hadoop and hive is proposed. The results showed that this Hadoop based system is more efficient and faster when compared with previously proposed solutions.

Research by, [5] discusses a methodology to migrate data relational data storage system like MySQL to from NoSQL databases, such as MongoDB and HBase using Hadoop MapReduce framework on a Hadoop Distributed File System (HDFS). The paper introduces Sqoop, a big data tool, that is used to achieve this migration. The migrated data was further analysed and moved into MongoDB using Hive.

[6] proposed implementation of the Spark framework for analyzing Twitter data. Spark was able to process the tasks within less time which shows the high efficiency of the framework. Spark uses storage such as HDFS for input-output operations. The high-speed performance of spark is contributed by the high rate of access of disk in a second and storing of current data in memory. These features encourage the use of Spark for data processing.

A research survey by, [7] on parallel processing systems explains the working of MapReduce architecture and its efficiency for data processing. It introduces Spark as a clusters computing framework, which is easy to use. It provides parallel operation features such as reduce, foreach and collect that facilities users to use functions to perform filter, map and reduce operations. Further, the author highlights Hive, a warehouse, which provides SQL like processing also known as Hive Query Language(HQL). Hive uses the directory in HDFS to store tables.

Similarly, [2] discusses the use of Apache Pig for data processing which is a high-level scripting language. The pig can access data using HDFS and perform execution parallelly. In contrast to SQL, there is no requirement for the schema in Pig and it can process both semi-structured and unstructured data.

## III. Methodology

### A. Dataset

The dataset used is downloaded from https://data. world/aprasla0922/online-retail#. It contains transaction details of an Online Retail store with over 541,909 rows. The columns for the dataset are incoiceNo, stockCode, description, quantity, invoiceDate, unitPrice, customerId and country.

### B. Data Pre-Processing

In this stage, cleansing of the dataset was performed using the R programming language in RStudio. The process involved the following steps.

1) Generating ID: To identify each row uniquely for analysis ID column was dynamically generated for each record in the table.

2) Splitting column: The invoice date column was split to form year, month, day and time columns.
3) Calculating Amount: A new column, 'amount' was created in which the total cost of each transaction was calculated stored using unitPrice and quantity columns.
4) Calculating Quarter: A new column, 'quarter' was created which describes the Quarter of Year in which a specific transaction was done.

After the completing of the above pre-processing tasks, the dataset is stored in a specified location in txt format for further analysis which is discussed in the next section.
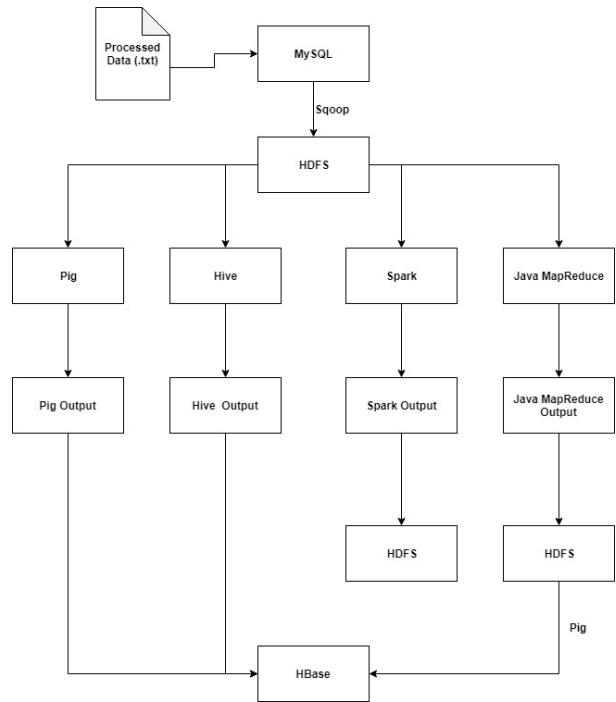
### C. Process Flow Approach



Fig. 1. Process Flow

As shown in 1, the processed data is stored in a MySQL database. Further, the data is moved to HDFS environment using Sqoop, in-order to load into Pig, Hive, Spark and Java MapReduce Frameworks. Each of the four frameworks, execute assigned MapReduce tasks (queries). The output generated by Pig script and Hive is directly stored into the HBase database. Whereas, outputs of Java MapReduce is first stored in HDFS before moving to HBase. To store Java MapReduce outputs into HBase, pig script has been used. The output of Spark SQL is stored in HDFS.

All the above tasks including cleaning of data using R script is automated except for Spark SQL, which is discussed in Section V. To automate the process, a shell script(.sh) has been used.

*1) Pig MapReduce:* :
The data for Pig queries were fetched from HDFS and the output were stored in HBase table. Figure 2
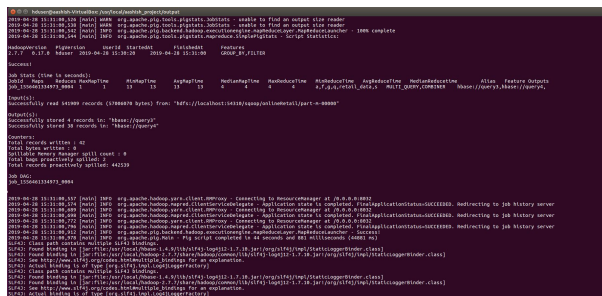
Fig. 2. Pig Script Completed



Fig. 5. Spark SQL in Jupyuter Notebook

shows successfully executed pig script with outputs stored in HBase.
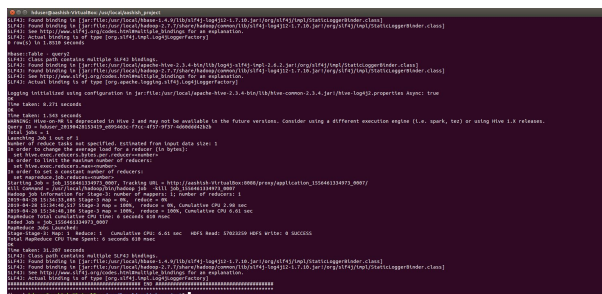
*2) Hive:* :



Fig. 3. Hive Script Completed

Similar to Pig, data for Hive was fetched from HDFS and inserted into Hive table. Figure 5 is a snapshot from the automated script showing executed Hive scripts.
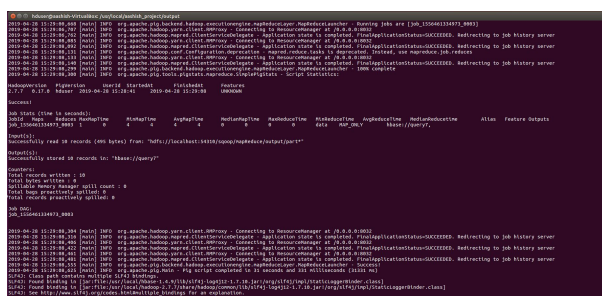
*3) Java MapReduce:* :



Fig. 4. Java MapReduce Completed

For Java MapReduce, Java code was complied into JAR file which was integrated into shell script for execution. Figure 5 shows completed Hive script.

*4) Spark:* :

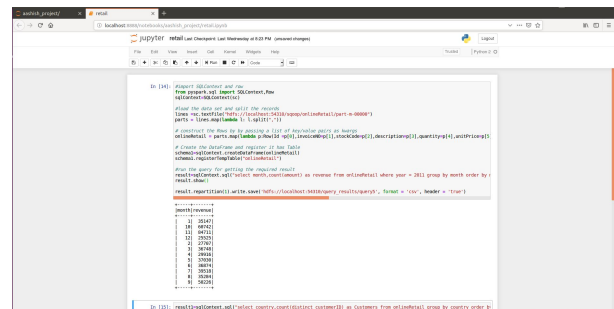For Spark SQL, Python code was used and the output was stored into HDFS.

*D. Technologies Used and Justification for choosing them*

The technologies used for the research are as mentioned below.

**Apache Hadoop**: Hadoop is an open-source framework that supports distributed computing. It is a Java implementation which runs on the MapReduce programming model and uses the Hadoop Distributed File System (HDFS)for data storage [4]. It ensures fault tolerance and eliminates data loss due to replication on server [2]. These features enable Hadoop to process data efficiently in less time and so it has been considered for this research.

**Hadoop MapReduce**: A programming model based on Java implementation for distributed computing. MapReduce algorithm involves map and reduces tasks. The map function takes data as input and split it before passing its output to the reducer. The research involves implementation of the MapReduce framework along with Pig, Hive and Spark.

**Apache Pig**: It is an open source framework that provides Pig Latin, a high-level programming language. It supports parallel processing of MapReduce tasks, hence reducing complexity [2]. The research uses pig script for MapReduce operation.

**Apache Hive**: A data warehouse that simplifies the use of Hadoop. It stores data in a structured format which is easy to understand from the user perspective. It provides SQL like query processing known as Hive query language (HQL) which provides accessibility and manipulation of stored data in HDFS or HBase. This makes Hive suitable for business applications [2].

**Apache Spark**: It is a distributed processing framework, that works on the in-memory system for fast processing. It is known for its high performance, easy to use and flexibility with efficiency in handling large datasets. Also, it supports application development in languages like python and java using Hadoop based storage system. One important feature

of Spark is that it provides Spark SQL, which allows users to write SQL commands to derive information from external data [2]. In this research, we use Spark SQL for query processing using the python programming language.

**Sqoop**: Sqoop is a big data tool that is used to transfer mass data from MySQL to HDFS. It uses MapReduce for import na dexport of data [5]. Being a fast and reliable tool for data transfer, we have considered the use of Scoop for data movement in our research.

**MySQL**: A Relational Database Management System to store structured data, known for its reliability, easy to use features, flexibility and high performance. We have used this database as initial storage, which acts as an input for other databases.

**Apache HBase**: HBase is an open source non-relational database which uses a distributed storage system for storing structured data. On top of Hadoop and HDFS, HBase provides capabilities of BigTable [1]. We have used HBase to store outputs generated by Pig, Hive, Spark and Java MapReduce.

**Java**: Java is a high-level object-oriented programming language. The MapReduce implementation is developed in java for its Map and reduce functionality. So, the research has used Java for MapReduce.

**R**: It is a programming language for statistical analysis and supports various operating system environments including UNIX [2]. We have used R data prepossessing and cleansing with RStudio software.

**Python**: Python is another powerful programming language which we have used with the Spark framework. PySpark is used to implement the Spark programming model using python code[3].

**PowerBI**: Power BI is a business analytics tool for data visualisation developed by Microsoft. It has been used in our research to represent various data outputs in graphical form.

**Tableau**: It is a data visualisation tool that is used for dynamic data representation. We have used this tool to represent data in Heat Map.

## IV. RESULTS

This section describes the results achieved for the objectives mentioned in Section I-A

**What are the top 20 most popular products among customers from the United Kingdom, sold in the first quarter of the year?**
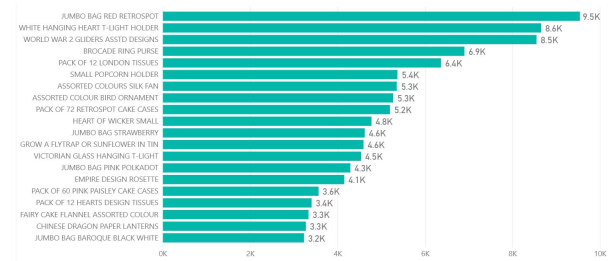


Fig. 6. Most popular Products

Figure 6 shows the top 20 most popular products purchased by customers from the United Kingdom. The units sold for these popular products range from 3.2K to 9.5K of which 'Jumbo Bag Red Retrospot' is the maximum.

**What is the average number of daily online orders received on an hourly basis?**
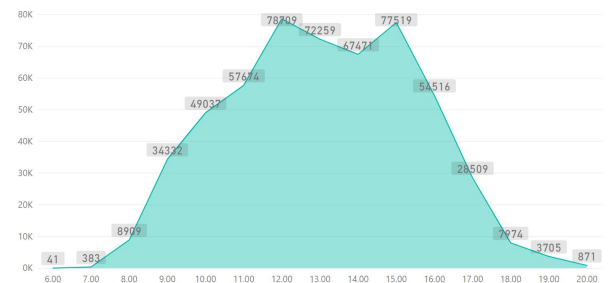


Fig. 7. Average number of orders per hour

As shown in figure 7, most of the orders are received between 12.00 to 15.00 daily. As users during these hours are most active, this duration can be used for promotional offers, special discounts and advertisement which shall reach mass users, improving the chances of business growth.

**How much was Net Profit on a Quarterly basis for the year 2011?**

Figure 8 describes the annual profit for the year on a quarterly basis. Quarter 4 accounted for the maximum profit with 32.96 percent of total net profit, followed by Quarter 3. This information gives an insight of increasing revenue with respect to time quarterly
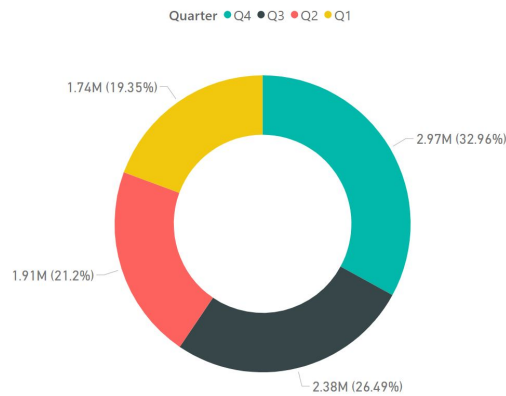
Fig. 8. Quarterly Net Profit

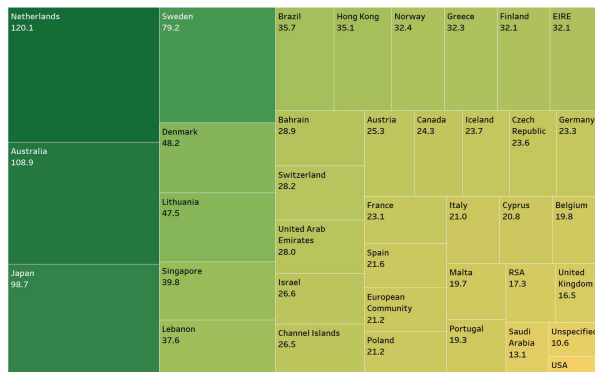## What is the average order amount for each country?



Fig. 9. Order Amount by Country

Figure 9, represents average order value for various countries. It can be seen that order value from the Netherlands and Australia is highest with 120 and 108 Euros. This shows high purchasing power of Customers from these countries, hence, products recommendation can be developed for such customers based on their purchasing habits which can increase the number of orders from these countries.

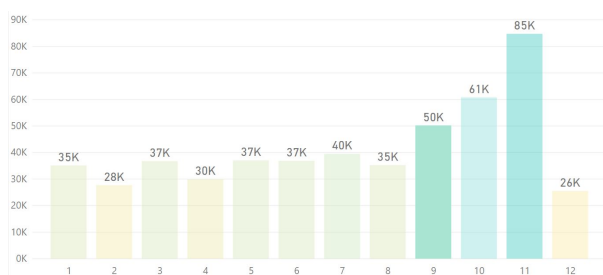## In which month the sales revenue reached a maximum?



Fig. 10. Monthly Sales Revenue

According to 10, sales revenue in the first half of the year was mostly levelled in the first half of the year. Whereas, in the second half revenue increased in the month of September, October and November. This shows that people tend to buy more in the second half of the year. This can be used as an opportunity for special offers and discounts.

## What is the total number of customers from each country?



Fig. 11. Number of Customer by Country

Figure 11 shows the number of customers from each county. Even though the order value for the United Kingdom is very less in comparison with other countries as seen in Figure 9, the number of customers is highest.

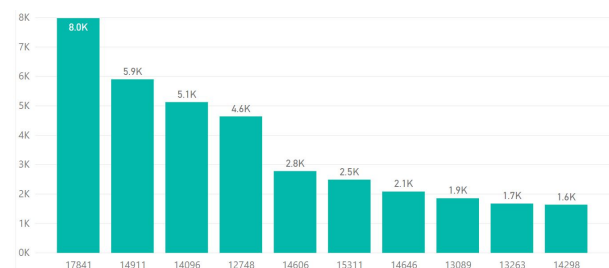## Which top 10 customer ID purchased the most?



Fig. 12. Top ten Customers with Maximum number of Purchases.

Figure 12 shows the top 10 customers who purchased maximum products. Hence, these customers tend to be high valued for the organisation.

### V. CHALLENGES AND LIMITATIONS

The development of the system involved many challenges. The most challenging part was integrating data cleaning with R script within automation process and saving Hive outputs into HBase within a single script.

Few limitations involve, compatibility issues, the Spark script could not be automated within the automation script. Unlike other outputs, the Spark outputs were stored into HDFS from where they were

moved to local directory. The data storage from Spark into HBase was not successful.

## VI. CONCLUSION AND FUTURE WORK

For analysing Big Data of online retail transactions, we developed a Hadoop MapReduce based system. The objective of the study was achieved as mentioned in Section I-A, to answer business queries. The system uses different MapReduce implementations using Pig, Hive, Java and Spark SQL with python which uses HDFS storage. The data was first imported into R for data cleansing before storing into MySQL database. Further, MySQL database was used as an input for other data processing frameworks. This demonstrate the capability of these platforms as a data processing frameworks.

The future work for this research should focus on including other useful parameters in dataset such as 'customer feedback' and 'product category' which can provide more meaningful insights. Also, the limitations in this research should be overcome by developing possible solutions for the issue mentioned in Section V.

## REFERENCES

[1] C. Jones and N. Livingstone, "Emerging implications of online retailing for real estate: Twenty-first century clicks and bricks," *Journal of Corporate Real Estate*, vol. 17, no. 3, pp. 226–239, 2015.

[2] A. Oussous, F.-Z. Benjelloun, A. A. Lahcen, and S. Belfkih, "Big data technologies: A survey," *Journal of King Saud University-Computer and Information Sciences*, vol. 30, no. 4, pp. 431–448, 2018.

[3] D. Gao, N. Wang, Z. He, and T. Jia, "The bullwhip effect in an online retail supply chain: a perspective of price-sensitive demand based on the price discount in e-commerce," *IEEE Transactions on Engineering Management*, vol. 64, no. 2, pp. 134–148, 2017.

[4] A. M. Hendawi, F. Alali, X. Wang, Y. Guan, T. Zhou, X. Liu, N. Basit, and J. A. Stankovic, "Hobbits: Hadoop and hive based internet traffic analysis," in *2016 IEEE International Conference on Big Data (Big Data)*. IEEE, 2016, pp. 2590–2599.

[5] P. M. Bante and K. Rajeswari, "Big data analytics using hadoop map reduce framework and data migration process," in *2017 International Conference on Computing, Communication, Control and Automation (ICCUBEA)*. IEEE, 2017, pp. 1–5.

[6] K. Aziz, D. Zaidouni, and M. Bellafkih, "Real-time data analysis using spark and hadoop," in *2018 4th International Conference on Optimization and Applications (ICOA)*. IEEE, 2018, pp. 1–6.

[7] Y. Zhang, T. Cao, S. Li, X. Tian, L. Yuan, H. Jia, and A. V. Vasilakos, "Parallel processing systems for big data: a survey," *Proceedings of the IEEE*, vol. 104, no. 11, pp. 2114–2136, 2016.