



Ministry of Communication  
and Information Technology



Next Technology Leaders

Intensive programs on Communication and Information  
Technology (3 Months)

Track  
Big Data Science

**Recommendation Application**

DEC 2018

# **Recommendation Application**

**BY**

Mohamed Sayed Abdelwahab

Mahmoud Abdullah Tawfik

Abdelrahman Mohmmmed Abdelmoniem

Mohammed Helal

**Under the supervision of**

**Dr. Doaa Hassan**

**Eng. Sayed A. Omar**

## **Acknowledgements**

We have taken efforts in this project. However, it would not have been possible without the kind support and help of many individuals and organizations. We would like to extend our sincere thanks to all of them.

We are highly indebted to our supervisors: Prof. Doaa Hassan and Eng. Sayed for their guidance and constant supervision as well as for providing necessary information regarding the project and also for their support in completing the project.

We would like to express our gratitude towards members of the National telecommunication Institute for their kind co-operation and providing us with necessary equipment for learning Data science Track.

Our thanks and appreciations also go to our colleagues in developing the project and people who have willingly helped us out with their abilities.

M.Aboalarbe, Mahmoud, Abdelrahman, Mohammed

# **Abstract**

In this project, we explore the technique of recommendation engine and its different types on Market basket to recommend items of the online Market to the most desirable customers. Using a dataset obtained from Instacart.com, we apply Association rule with Apriori Algorithm of about 7000000 rows. We compare performance time of the project when it run with data science and with Apache spark using databricks.com, then we make an API model and connect it with a pretty graphical user interface in an Android Application to be more friendly to the users.

## List of figures

<b>Figures</b>	<b>Number of pages</b>
Recommendation engine uses	10
Association Rule	12
Apriori algorithm	13
Logistic regression in Hadoop and Spark (Speed)	14
Apache Spark Uses	15
project Architecture	19
Dataset description	20
Highest department	21
Most sold item	22
Most reordered Item	23
Highest selling hour	24
import libraries	25
association rule	25
result	26
import libraries	26
association rule	27
result	27
Most freq item A	28
Lift for most freq items A & B	28
Top items with highest lift	29
API	30
GUI in Android APP	31

# Table of Contents

<b>Acknowledgement.....</b>	<b>3</b>
<b>Abstract.....</b>	<b>4</b>
<b>List of figures.....</b>	<b>5</b>
<b>Chapter 1 Introduction.....</b>	<b>7</b>
1.1 project main idea.....	8
1.2 problem Statement.....	8
1.3 proposed solution.....	8
<b>Chapter 2 Background and Theoretical Overview.....</b>	<b>9</b>
2.1 Scientific background.....	10
2.1.1 Recommendation Engines.....	11
2.1.2 Association rule.....	11
2.1.3 Apriori Algorithm.....	13
2.1.4 Apache Spark.....	14
2.1.5 Android.....	16
<b>Chapter 3 Project Implementation.....</b>	<b>17</b>
3.1 Used Tools.....	18
3.2 Project Archtecture.....	19
3.3 Dataset Description.....	20
3.4 Dataset Visualization .....	21
3.5 Model Implementation .....	25
3.5.1 IN Data Science Model.....	25
3.5.2 IN Big Data Model.....	26
3.6 Model Evaluation .....	28
3.7 API Implementation.....	30
3.8 project Demo (Recommendation Android App) .....	31
<b>Chapter 4 Conclusions.....</b>	<b>35</b>
<b>References.....</b>	<b>36</b>

# **Chapter one**

## **Introduction**

# **Chapter One**

## **Introduction**

### **1.1 project main idea**

The main idea of the project is to make a recommendation application to recommend grocery product for customers so the market could save its money spent advertising and customers save their time for searching for their desired products.

### **1.2 problem Statement**

Companies spent money on advertising and sometimes high portion of this money doesn't make an effect because these advertising go to the wrong customers that don't prefer these products or don't match their needs, so we want to make an application that would recommend products to the most desirable Customers that have high probability to buy these items.

### **1.3 proposed solution**

Our solution is to make a recommendation application that recommend the grocery items to the customers according to their needs, so company can recommend the best products for their customer that have high probability to buy these products.



## **Chapter Two**

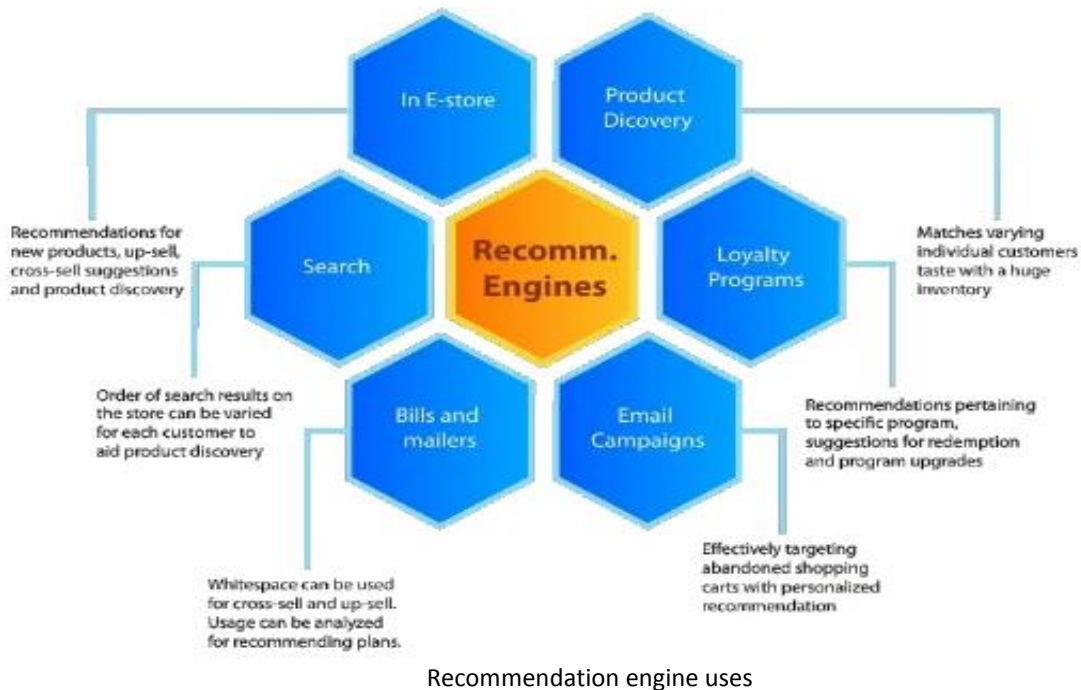
### **Background and Theoretical Overview**

# Chapter Two

## Background and Theoretical Overview

### 2.1 Scientific background

#### 2.1.1 Recommendation Engines



In today's world, every customer is faced with multiple choices. For example, If I'm looking for a book to read without any specific idea of what I want, there's a wide range of possibilities how my search might pan out. I might waste a lot of time browsing around on the internet and trawling through various sites hoping to strike gold. I might look for recommendations from other people.

But if there was a site or app which could recommend me books based on what I have read previously, that would be a massive help. Instead of wasting time on various sites, I could just log in and voila! 10 recommended books tailored to my taste.

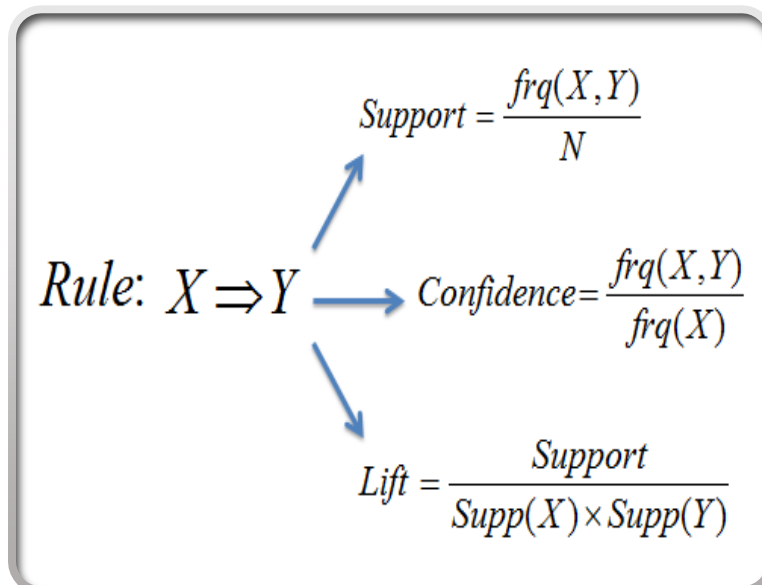
This is what recommendation engines do and their power is being harnessed by most businesses these days. From Amazon to Netflix, Google to Goodreads, recommendation engines are one of the most widely used applications of machine learning techniques. A recommendation engine filters the data using different algorithms and recommends the most relevant items to users. It first captures the past behavior of a customer and based on that, recommends products which the users might be likely to buy.

If a completely new user visits an e-commerce site, that site will not have any past history of that user. So how does the site go about recommending products to the user in such a scenario? One possible solution could be to recommend the best-selling products, i.e. the products which are high in demand. Another possible solution could be to recommend the products which would bring the maximum profit to the business.

If we can recommend a few items to a customer based on their needs and interests, it will create a positive impact on the user experience and lead to frequent visits. Hence, businesses nowadays are building smart and intelligent recommendation engines by studying the past behavior of their users.

### 2.1.2 Association rule

Association Rules find all sets of items (itemsets) that have support greater than the minimum support and then using the large itemsets to generate the desired rules that have confidence greater than the minimum confidence. The lift of a rule is the ratio of the observed support to that expected if X and Y were independent. A typical and widely used example of association rules application is market basket analysis



Association Rule

#### Support (S):

Support is an indication of how frequently the items appear in the database.

#### Confidence (C):

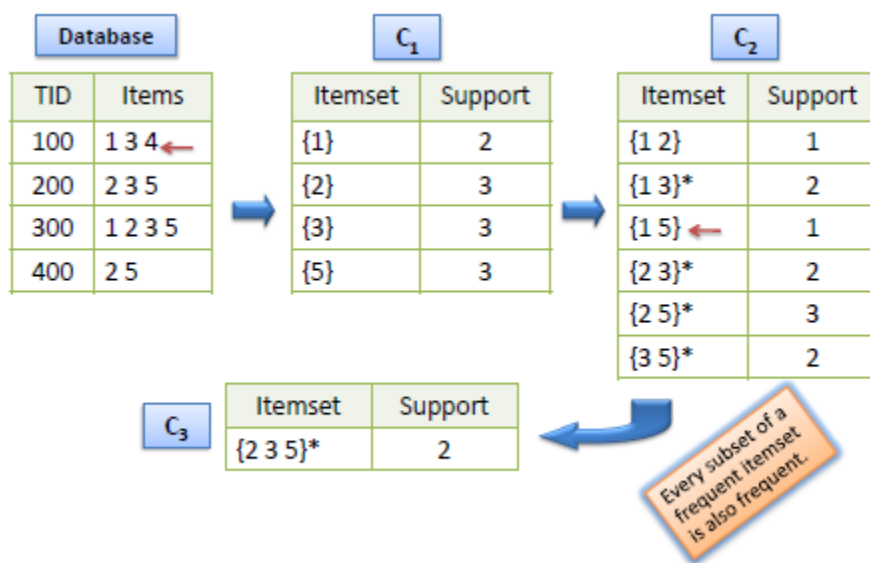
indicates the number of times the if/then statements have been found to be true.

#### Lift:

Lift is the ratio of the observed support to that expected if X and Y were independent.

### 2.1.3 Apriori Algorithm

1. Candidate itemsets are generated using only the large itemsets of the previous pass without considering the transactions in the database.
2. The large itemset of the previous pass is joined with itself to generate all itemsets whose size is higher by 1.
3. Each generated itemset that has a subset which is not large is deleted. The remaining itemsets are.



Apriori algorithm

The Apriori algorithm takes advantage of the fact that any subset of a frequent itemset is also a frequent itemset. The algorithm can therefore, reduce the number of candidates being considered by only exploring the itemsets whose support count is greater than the minimum support count. All infrequent itemsets can be pruned if it has an infrequent subset.

## 2.1.4 Apache Spark

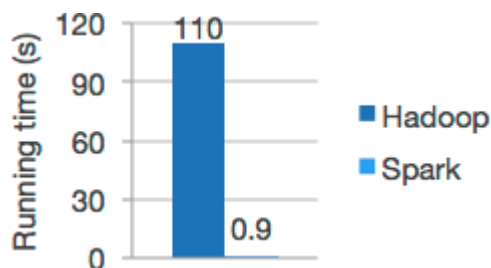
Apache Spark is an open-source distributed general-purpose cluster-computing framework. Originally developed at the University of California, Berkeley's AMPLab, the Spark codebase was later donated to the Apache Software Foundation, which has maintained it since. Spark provides an interface for programming entire clusters with implicit data parallelism and fault tolerance.

### Advantages of Apache Spark

#### Speed

Run workloads 100x faster.

Apache Spark achieves high performance for both batch and streaming data, using a state-of-the-art DAG scheduler, a query optimizer, and a physical execution engine.



Logistic regression in Hadoop and Spark (Speed)

#### Ease of Use

Write applications quickly in Java, Scala, Python, R, and SQL.

Spark offers over 80 high-level operators that make it easy to build parallel apps. And you can use it *interactively* from the Scala, Python, R, and SQL shells.

#### Generality

Combine SQL, streaming, and complex analytics.

Spark powers a stack of libraries including SQL and DataFrames, Mllib for machine learning, GraphX, and Spark Streaming. You can combine these libraries seamlessly in the same application.

## Runs Everywhere

Spark runs on Hadoop, Apache Mesos, Kubernetes, standalone, or in the cloud. It can access diverse data sources.

You can run Spark using its standalone cluster mode, on EC2, on Hadoop YARN, on Mesos, or on Kubernetes. Access data in HDFS, Alluxio, Apache Cassandra, Apache Hbase, Apache Hive, and hundreds of other data sources.



Apache Spark Uses

### 2.1.5 Android (operating system)



Android is a mobile operating system developed by Google, based on a modified version of the Linux kernel and other open source software and designed primarily for touchscreen mobile devices such as smartphones and tablets. In addition, Google has further developed Android TV for televisions, Android Auto for cars, and Wear OS for wrist watches, each with a specialized user interface. Variants of Android are also used on game consoles, digital cameras, PCs and other electronics.

Initially developed by Android Inc., which Google bought in 2005, Android was unveiled in 2007, with the first commercial Android device launched in September 2008. The operating system has since gone through multiple major releases, with the current version being 9 “Pie”, released in August 2018. The core Android source code is known as Android Open Source Project (AOSP), and is primarily licensed under the Apache License.

Android is also associated with a suite of proprietary software developed by Google, called Google Mobile Services (GMS) that very frequently comes pre-installed in devices, which usually includes the Google Chrome web browser and Google Search and always includes core apps for services such as Gmail, as well as the application store and digital distribution platform Google Play, and associated development platform. These apps are licensed by manufacturers of Android devices certified under standards imposed by Google, but AOSP has been used as the basis of competing Android ecosystems, such as Amazon.com’s Fire OS, which use their own equivalents to GMS.



## **Chapter Three**

### **Project Implementation**

## Chapter Three

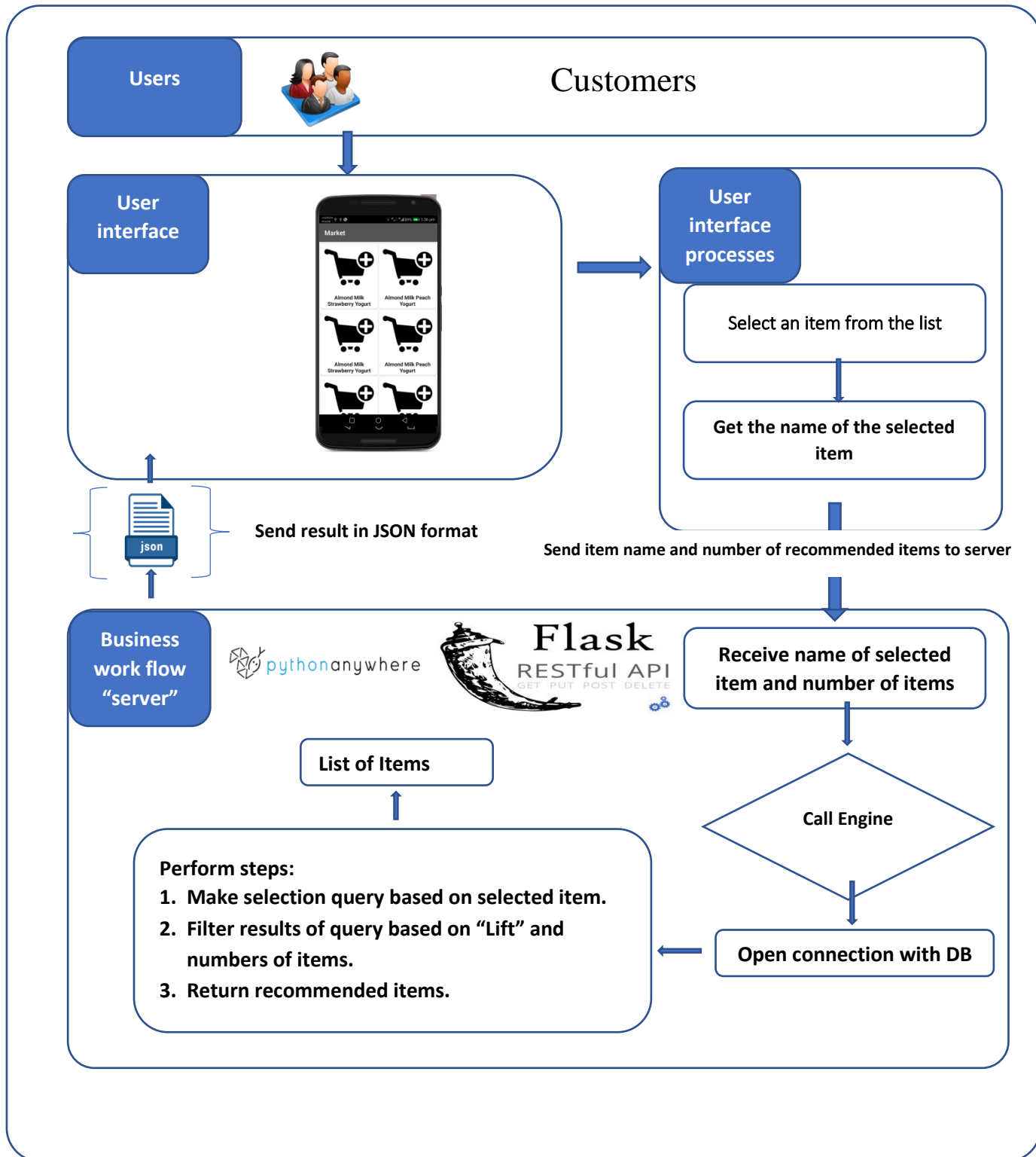
### Project Implementation

#### 3.1 Used Tools

Our tools that was used in this project:

- Python v3 programming language
- Android studio
- Python Flask for API on PythonAnywhere
- Jupyter Notebook
- Apache spark on Databricks
- Tableau software
- API.

## 3.2 Project Architecture



project Architecture

### 3.3 Data set description

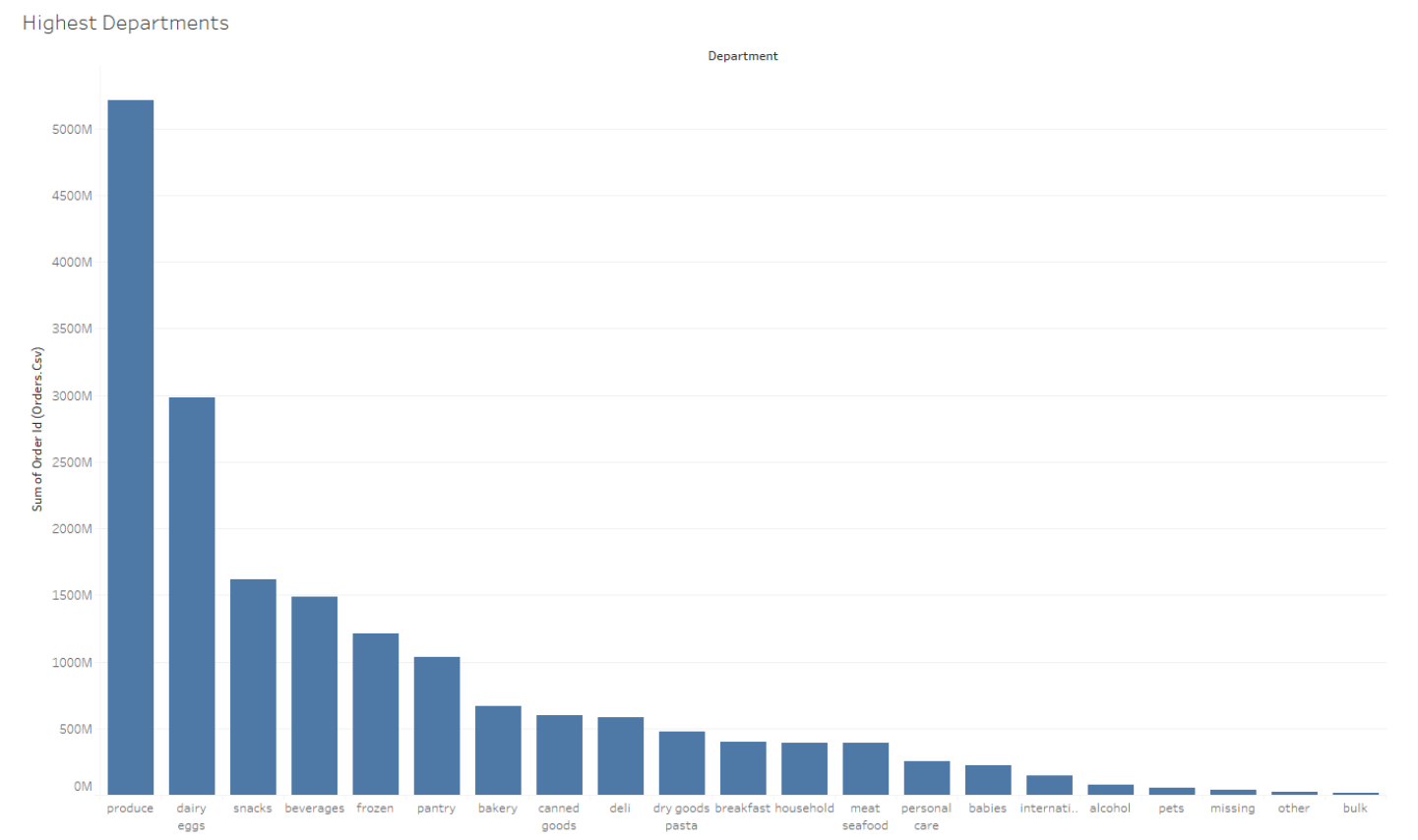
Our dataset consists of 32434489 rows and four features

order_id	product_id	add_to_cart_order	reordered
2	33120	1	1
2	28985	2	1
2	9327	3	0
2	45918	4	1
2	30035	5	0
2	17794	6	1
2	40141	7	1
2	1819	8	1
2	43668	9	0
3	33754	1	1
3	24838	2	1
3	17704	3	1
3	21903	4	1
3	17668	5	1
3	46667	6	1
3	17461	7	1
3	32665	8	1
4	46842	1	0
4	26434	2	1
4	39758	3	1
4	27761	4	1

The URL of dataset: <https://www.instacart.com/datasets/grocery-shopping-2017>

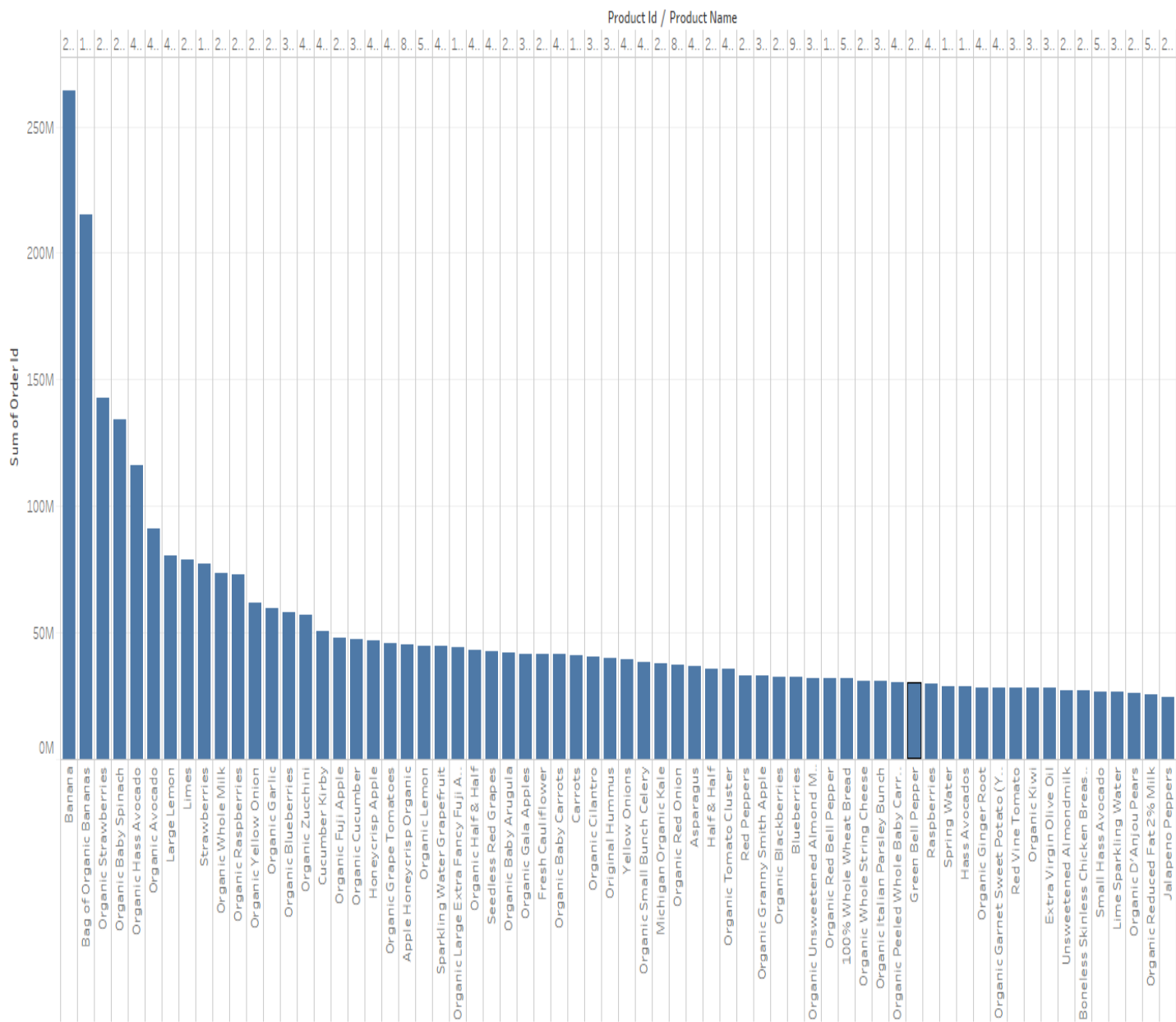
### 3.4 Data set visualization:

Using **Tableau**, we have reached some visualizations:



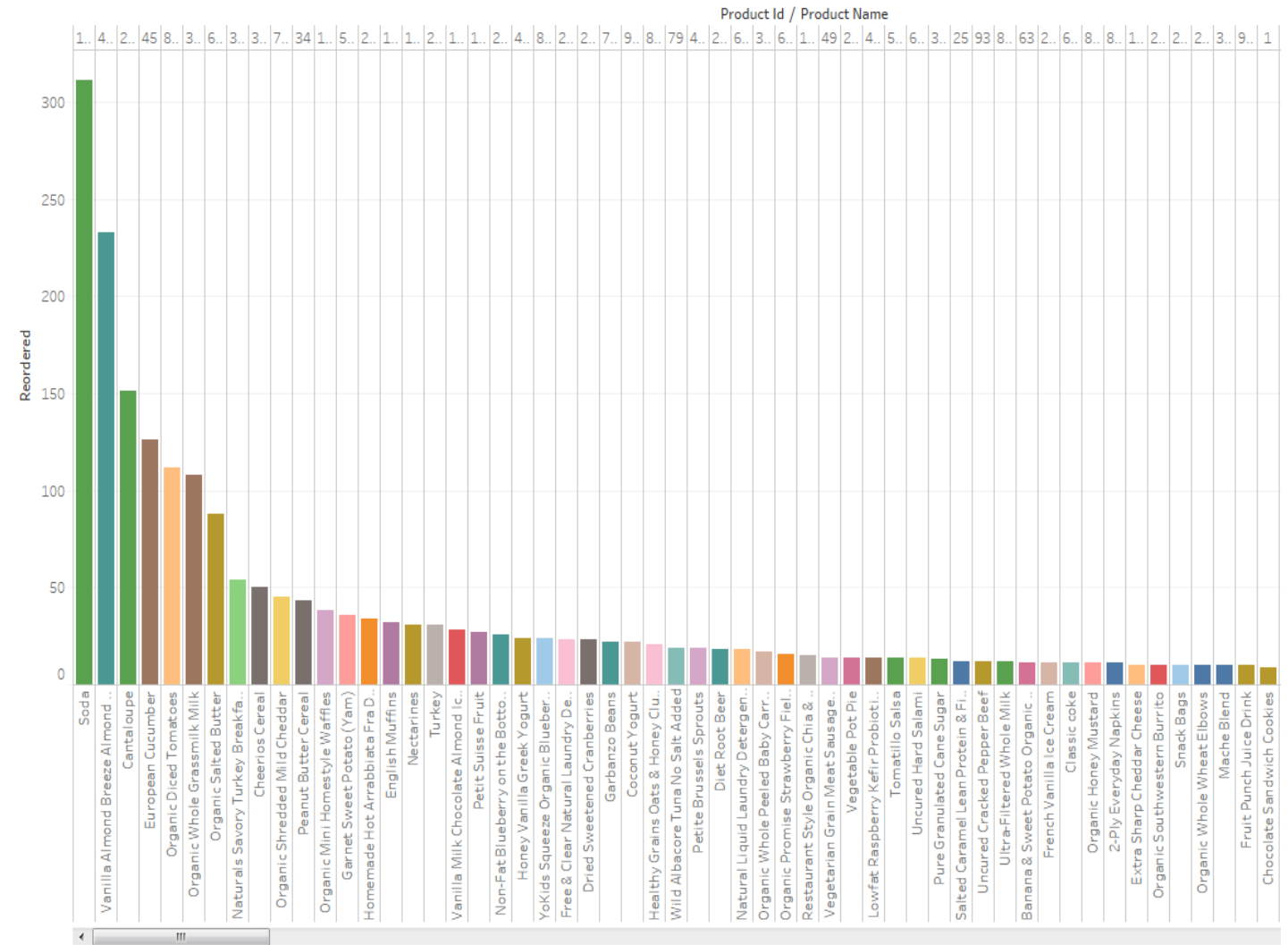
Highest department

Most selling items

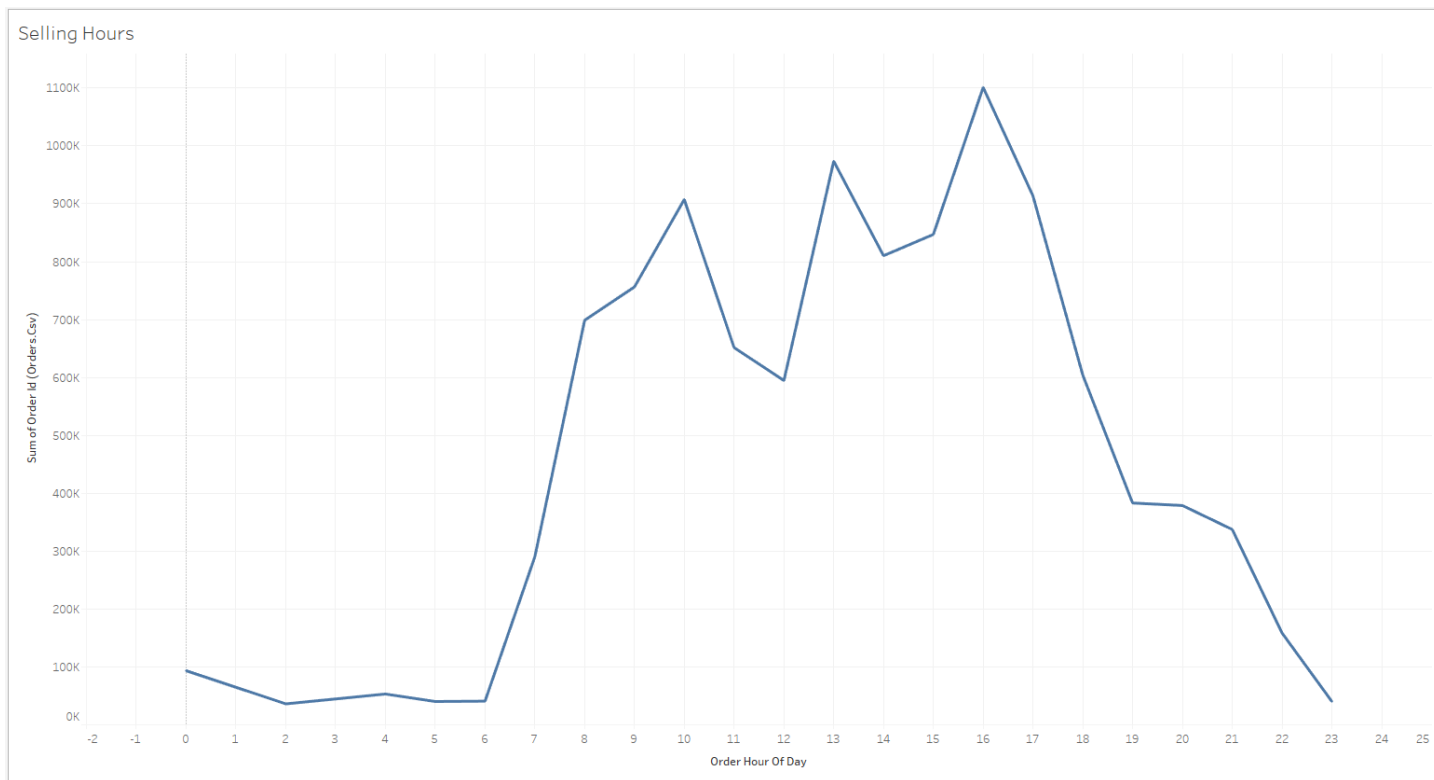


Most sold items

## Most reordered Item



Most reordered Item



Highest selling hour



## 3.5 Model Implementation

### 3.5.1 IN Data Science Model

- We have worked on dataset which consist of 7000000 rows and two features to be able to analyze and process this dataset using Jupyter Notebook with python

We have used association rule for making this recommendation system with Apriori Algorithm

First:

we import libraries in our project

```
import pandas as pd
import numpy as np
import sys
from itertools import combinations, groupby
from collections import Counter
from IPython.display import display
```

import libraries

Second:

The main function of the system which uses association rule

```
def association_rules(order_item, min_support):

    print("Starting order_item: {:22d}".format(len(order_item)))

    # Calculate item frequency and support
    item_stats = freq(order_item).to_frame("freq")
    item_stats['support'] = item_stats['freq'] / order_count(order_item) * 100

    # Filter from order_item items below min support
    qualifying_items = item_stats[item_stats['support'] >= min_support].index
    order_item = order_item[order_item.isin(qualifying_items)]

    print("Items with support >= {}: {:15d}".format(min_support, len(qualifying_items)))
    print("Remaining order_item: {:21d}".format(len(order_item)))
```

association rule

After using Association rule, we reach this **result**:

```
Starting order_item:          7000000
Items with support >= 0.04:    4123
Remaining order_item:         5494313
Remaining orders with 2+ items: 627718
Remaining order_item:         5443208
Item pairs:                   6422989
Item pairs with support >= 0.04: 6530

Wall time: 10min 29s
```

result

### 3.5.2 IN Big Data Model

We have worked on dataset which consist of 7000000 rows and two features to be able to analyze and process this dataset using Apache Spark on Databricks

**First:**

we import libraries in our project

```
Cmd 1
1 import math
2 import numpy
3 from pyspark.mllib.linalg import Vectors
4 from pyspark.mllib.regression import LabeledPoint
5 from pyspark.mllib.stat import Statistics
6 from pyspark.sql import SQLContext
7 from numpy import array
8 from pyspark.sql import Row
```

## Second:

The main function of the system which uses association rule

```
1 def association_rules(order_item, min_support):
2
3     print("Starting order_item: {:22d}".format(len(order_item)))
4
5
6     # Calculate item frequency and support
7     item_stats = freq(order_item).to_frame("freq")
8     item_stats['support'] = item_stats['freq'] / order_count(order_item) * 100
9
10
11     # Filter from order_item items below min support
12     qualifying_items = item_stats[item_stats['support'] >= min_support].index
13     order_item = order_item[order_item.isin(qualifying_items)]
14
15     print("Items with support >= {:}: {:15d}".format(min_support, len(qualifying_items)))
16     print("Remaining order_item: {:21d}".format(len(order_item)))
17
```

association rule

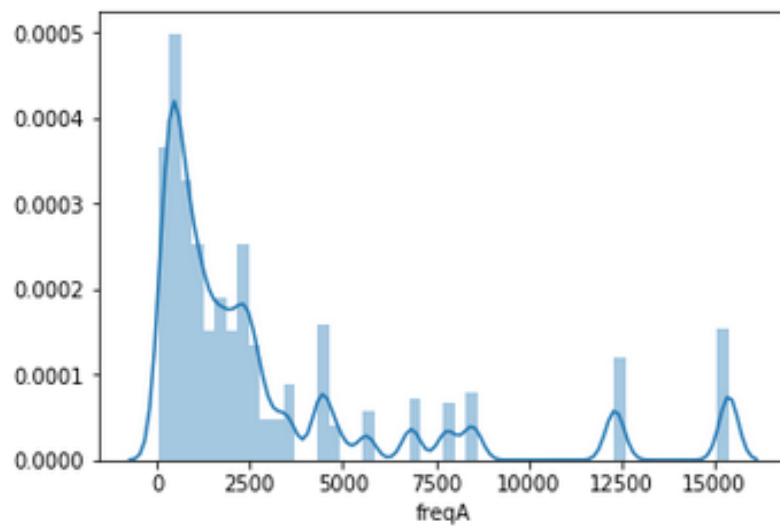
After using Association rule with Apache Spark on Data bricks we reach this **result**:

```
Starting order_item:          7000000
Items with support >= 0.04:    4123
Remaining order_item:        5494313
Remaining orders with 2+ items: 627718
Remaining order_item:        5443208
Item pairs:                  6422989
Item pairs with support >= 0.04: 6530

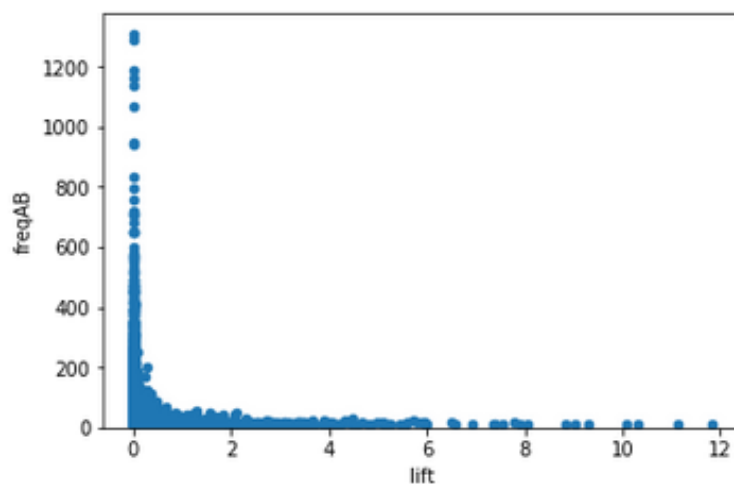
CPU times: user 1min 32s, sys: 6.64 s, total: 1min 38s
Wall time: 1min 36s
```

result

### 3.6 Model Evaluation

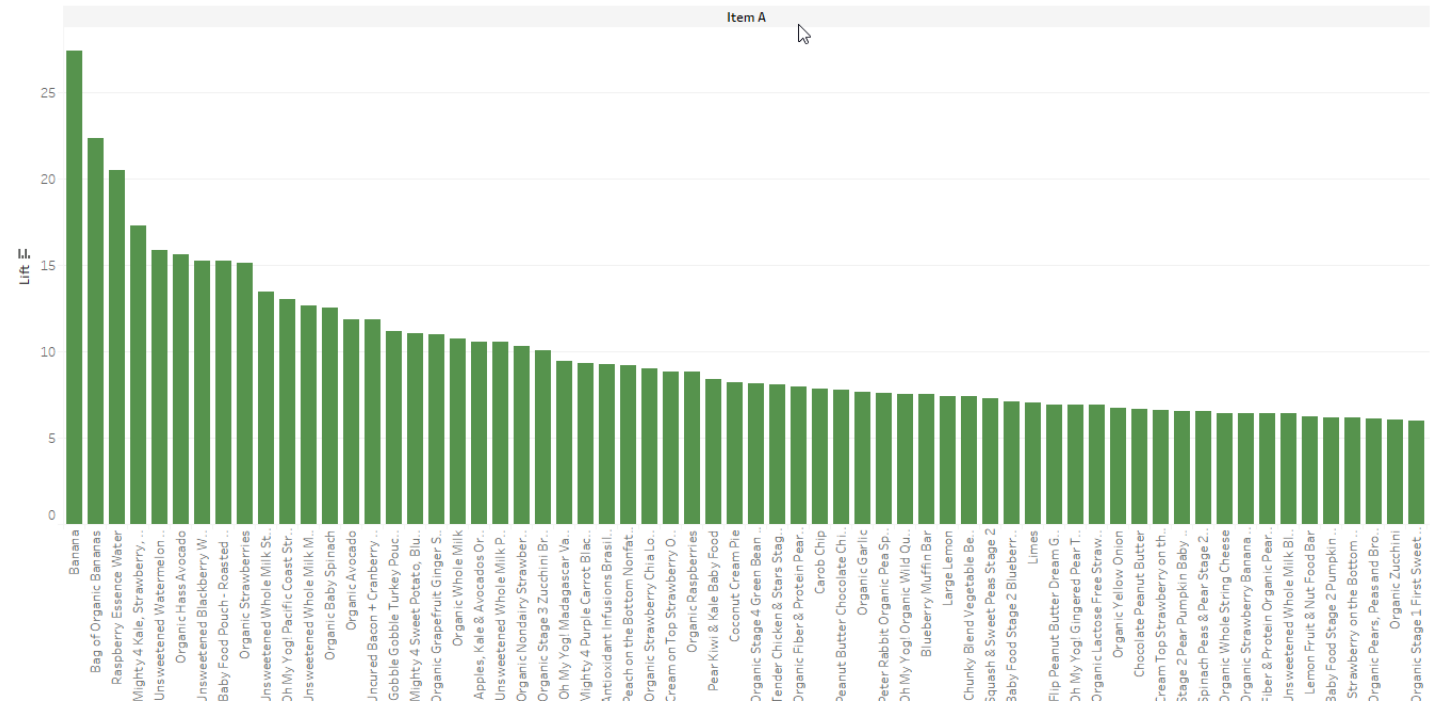


Most freq item A



Lift for most frequent item A with B

Top item has the Highest Lift



Top item has the highest lift

### 3.7 API Implementation:

We use **PythonAnywhere** is a rule-based machine learning method for discovering interesting relations between variables in large databases. It is intended to identify strong rules discovered in databases using some measures of interestingness. This rule-based approach also generates new rules as it analyzes more data. The ultimate goal, assuming a large enough dataset, is to help a machine mimic the human brain's feature extraction and abstract association capabilities from new uncategorized data.

```
mysql = MySQL()

class Database:
    def __init__(self):
        app.config['MYSQL_DATABASE_USER'] = 'recommendation'
        app.config['MYSQL_DATABASE_PASSWORD'] = '01005500593Mm'
        app.config['MYSQL_DATABASE_DB'] = 'recommendation$recommendation'
        app.config['MYSQL_DATABASE_HOST'] = 'recommendation.mysql.pythonanywhere-services.com'
        mysql.init_app(app)
        self.conn = mysql.connect()
        self.cursor = self.conn.cursor()

    def list_items(self):
        self.cursor.execute("SELECT itemA FROM rules LIMIT 30")
        result = self.cursor.fetchall()
        data = {"data":result}
        return jsonify(data)

    def get_recommendation_items(self,selected_item,no_items):
        self.cursor.execute("SELECT itemB FROM rules WHERE itemA = '"+selected_item+"' ORDER BY lift DESC LIMIT "+no_items+"")
        result = self.cursor.fetchall()
        data = {"data":result}
        return jsonify(data)
```

API

### 3.8 project Demo (Recommendation Android App)

We Use Graphical User Interface (GUI) using Android Application with friendly interface to the customers showing all the items and make him to choose item and then Application would recommend to the user the item based on his choice.



GUI in Android APP

### First:

We make Android Application with friendly Graphical user interface that is item-based recommendation engine in background and recommend items to the user.



GUI in Android APP

### Second:

We measure run-time and see difference with Apache Spark on Databricks, python on local machine and Cloudera VM and these are the results:

	Processing Time	No. of data items
Python	10 min 29 Sec	7000000
Apache Spark	1 min 36 Sec	7000000
Apache Spark	2 min 11 Sec	10000000
Apache Spark	4 min 2 Sec	20000000
Cloudera VM	1 min 35 Sec	7000000
Cloudera VM	Out of memory	10000000



We have run the algorithm on local machine of the following properties:

- Operating System: Windows 10 Pro 64-bit (10.0).
- Processor: Intel(R) Core(TM) i5-2520M CPU @ 2.50GHz (4 CPUs), 2.5GHz.
- Memory: 6G RAM.
- Display Memory: 1664 MB.
- Shared Memory: 1632 MB.

# **Chapter Four**

## **Conclusions**

## **Chapter Four**

### **Conclusions**

This section summarizes the conclusions of the project. It discusses both the technical and general findings of the project, this project succeeded in creating an android application to recommend Market items using association rule with Apriori Algorithm using dataset with from Instacart.com. We worked on dataset that consist of 32434489 rows and we worked on nearly 7000000 due to limitation of hardware resources of machine, we compare performance time of the project when it run with data science and with Apache spark using databricks.com, then we make an API model and connect it with a pretty graphical user interface in an Android Application to be more friendly to the users.

## References:

- [https://www.saedsayad.com/association\\_rules.htm](https://www.saedsayad.com/association_rules.htm)
- [https://en.wikipedia.org/wiki/Association\\_rule\\_learning](https://en.wikipedia.org/wiki/Association_rule_learning)
- Data Mining concepts and techniques by Jiawei Han
- <https://en.wikipedia.org/wiki/PythonAnywhere>
- [https://en.wikipedia.org/wiki/Machine\\_learning](https://en.wikipedia.org/wiki/Machine_learning)
- <https://www.adobe.com/experience-cloud/topics/recommendation-engine.html>
- <https://towardsdatascience.com/a-gentle-introduction-on-market-basket-analysis-association-rules-fa4b986a40ce>
- Github Repo Link: <https://github.com/Aboalarbe/Grocery-Recommendation-System>