# Global Superstores Data Exploration and Analysis for website Recommender System

Malungisa Mndzebele, DS-490-2021SP: Senior Project, Instructor: Dr. Shafqat Ali Shad

## Introduction

**Motivation:**

The main problem to be considered in this project is; Using grocery store customer and transaction data, can we make a reliable recommendation system for future customers? I thought it would be important to do my senior project focusing on this question because I think since all businesses depend on some form of recommender system to ensure efficiency, this would be a great skill to have. Hopefully by the end of this project I should have a good basic understanding of how recommender system like that of amazon or Netflix works.

Some insights from this project should also include the relationship between products bought by customers. I plan to make an interactive page where a customer can add a product to their cart and based on that the system should recommend other products to buy based on customer preference and other customers' activities.

**Data:**

Word data - https://www.kaggle.com/paultimothymooney/latitude-and-longitude-for-every-country-and-state. The Global Superstore data can be found at: https://data.world/tableauhelp/superstore-data-sets. The data contains 24 columns and 51 290 rows.

The data has 17 415 unique customer ID so this means the some customers have bought more that one product whether at the same time or at different times. There is transactional data but there is no specific column for preference. This can be solved by assuming that the product a customer buys is what should be in their preference or add a product to someone's preference if they buy more than one unit of buy if multiple times.





## METHODOLOGY

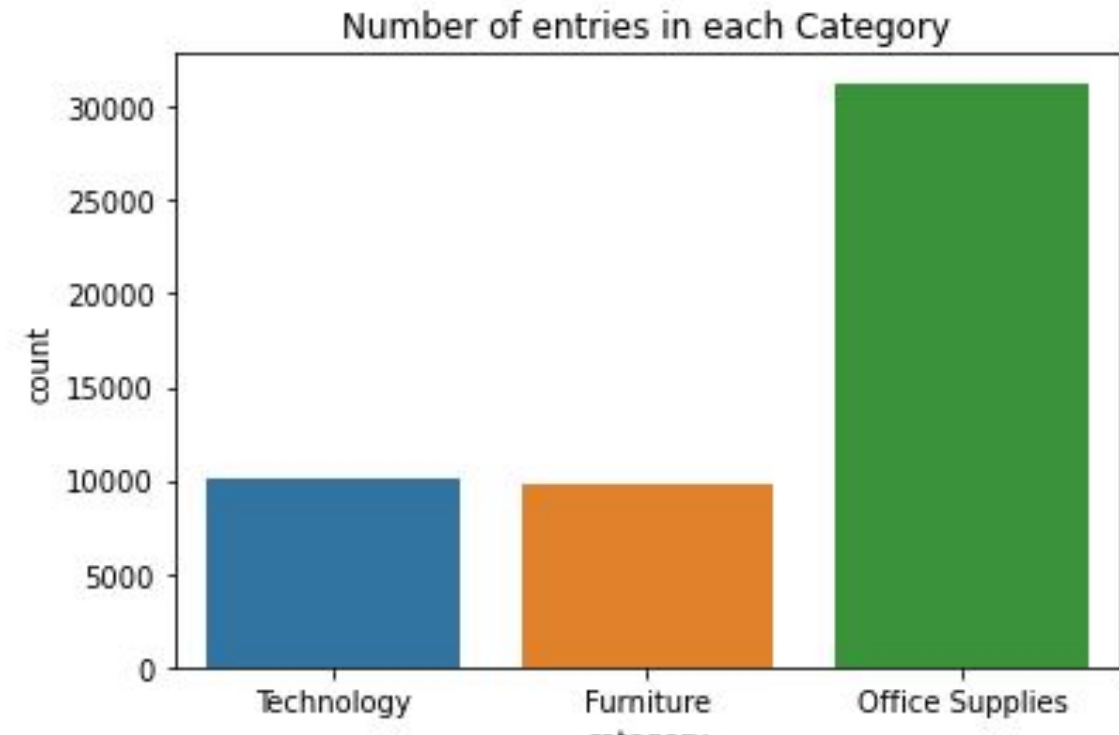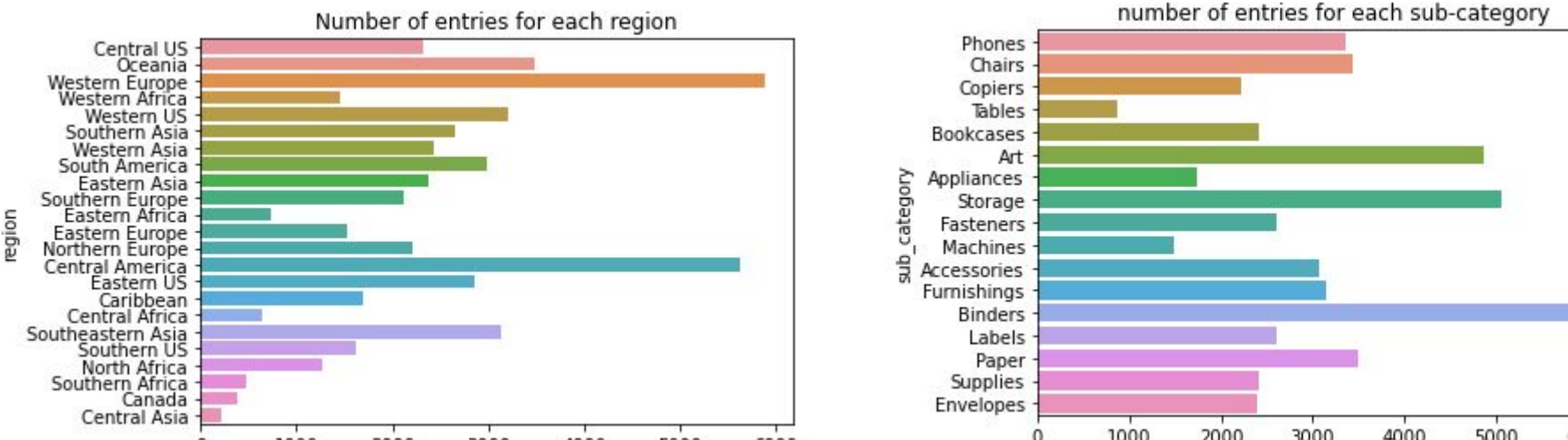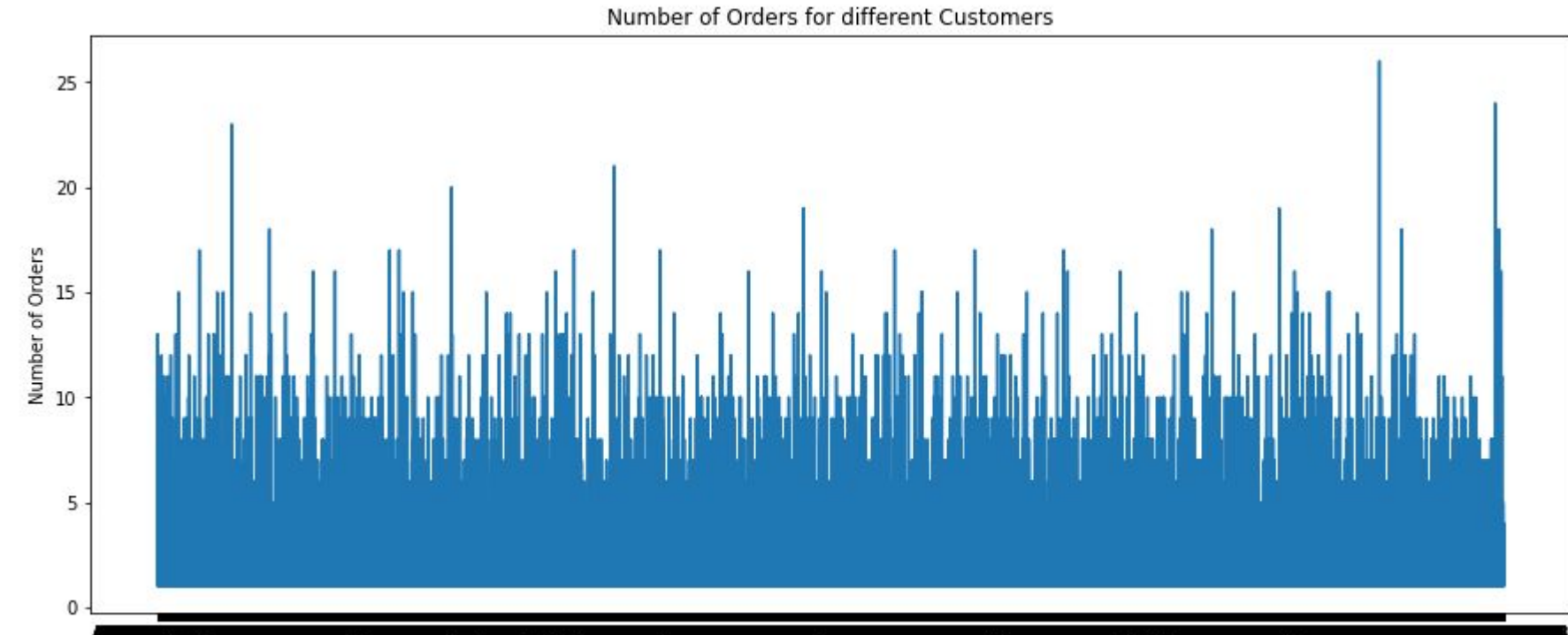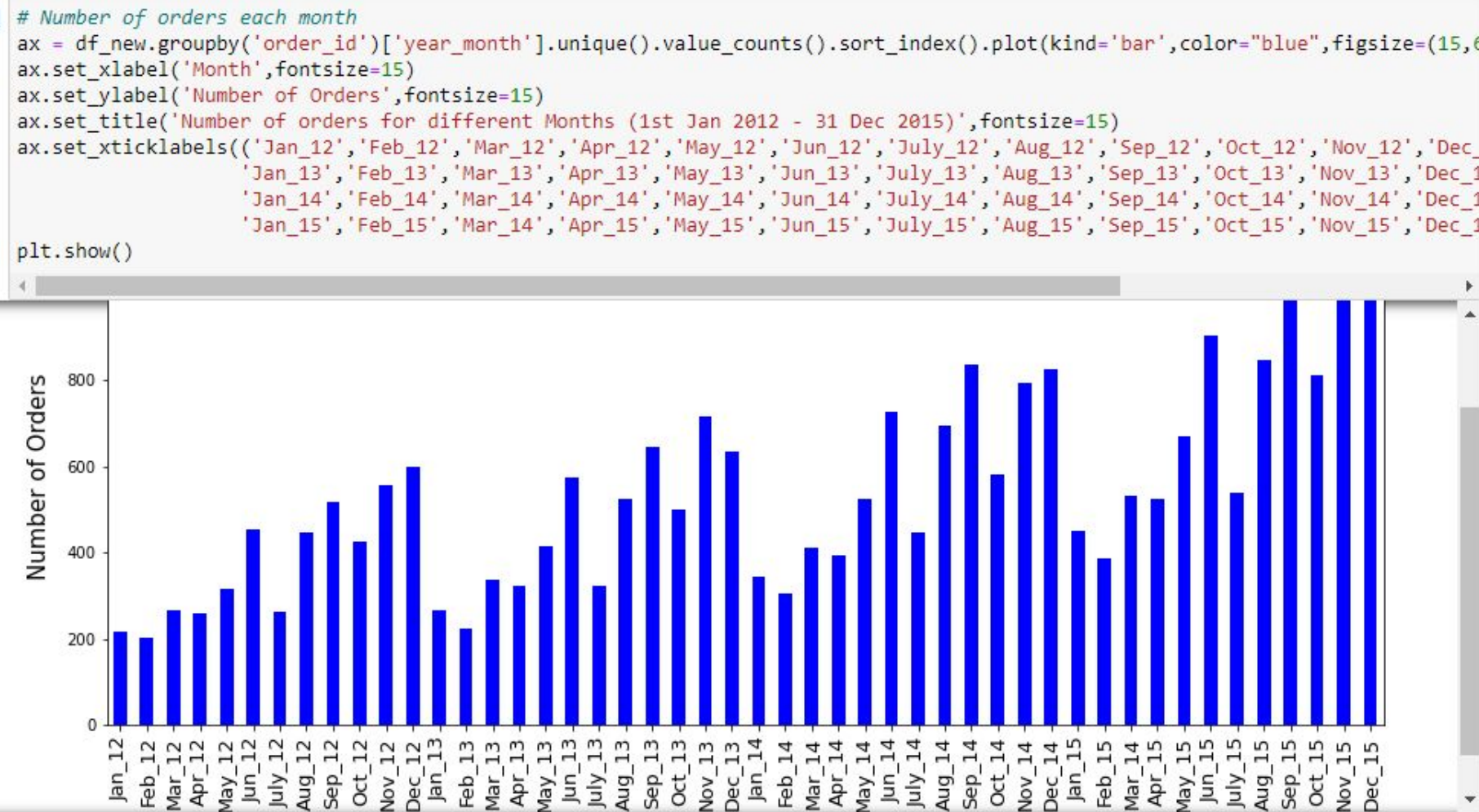Firstly I did some data Exploration and visualization in python .

Then model building in python: 1. Apriori algorithm for category and product recommendation. 2. K-means for market segmentation

Then Build a dashboard in Tableau to present Customer purchasing data, sales, profit and locational data.

Then implement a model on R for the recommender system.

## Data Exploration (python)









## MODEL Building (python)

**1. a. Apriori Algorithm for sub-category**



**1. b. Apriori Algorithm for each product sold**



**Tableau Dashboards:**



## Implementation of model in R