Detailed Project Report

# Customer Segmentation using k-prototype algorithm

Revision Number: 1.3

Last date of revision: 23/05/2022

**Rajdeep Mondal**

**Arpan Das**

# Abstract

This project aims to analyze [E-Commerce data](#) that list purchases made by nearly 4000 customers from December 2010 to December 2021. Based on this database we performed Exploratory Data Analysis with Statistical Methods for gaining data-driven insights with machine learning. Here we used Unsupervised techniques with Python for grouping the customers by their behavioral patterns.

## Dataset Description

This data contains 8 columns —

1. **InvoiceNo**: This is the Invoice number. There are 25,900 unique invoice data. It is a six-digit integral number uniquely assigned to each transaction. If this code starts with the letter 'C', it indicates a cancellation.
2. **StockCode**: This is the Product (item) code. There are 4,070 unique StockCode values. It is a five-digit integral number uniquely assigned to each distinct product. For some data, it contains special code like — `D`, `POST`, `M`, `C2`, `CRUK`, `Discount`, `POSTAGE`, `Manual`, `CARRIAGE`, `CRUK`, `Commission`.
3. **Description**: This describes the product, ie — Product Name. There are 4224 unique descriptions.
4. **Quantity**: This represents the quantities of each product (item) per transaction. It is a Numeric column.
5. **InvoiceDate**: This displays the Invoice Date and time which was generated when each transaction was completed. It is a Numeric column.
6. **UnitPrice**: This represents the Unit price of each product. It is a Numeric column.

7. **CustomerID**: This represents the unique Customer number. It is a five-digit integral number uniquely assigned to each customer.
8. **Country**: This represents the Country name where each customer resides.

```
1 <class 'pandas.core.frame.DataFrame'>
2 RangeIndex: 541909 entries, 0 to 541908
3 Data columns (total 8 columns):
4  #   Column       Non-Null Count   Dtype
5 ---  ------       --------------   -----
6  0   InvoiceNo    541909 non-null  object
7  1   StockCode    541909 non-null  object
8  2   Description  540455 non-null  object
9  3   Quantity     541909 non-null  int64
10 4   InvoiceDate  541909 non-null  object
11 5   UnitPrice    541909 non-null  float64
12 6   CustomerID   406829 non-null  float64
13 7   Country      541909 non-null  object
14 dtypes: float64(2), int64(1), object(5)
15 memory usage: 33.1+ MB
```
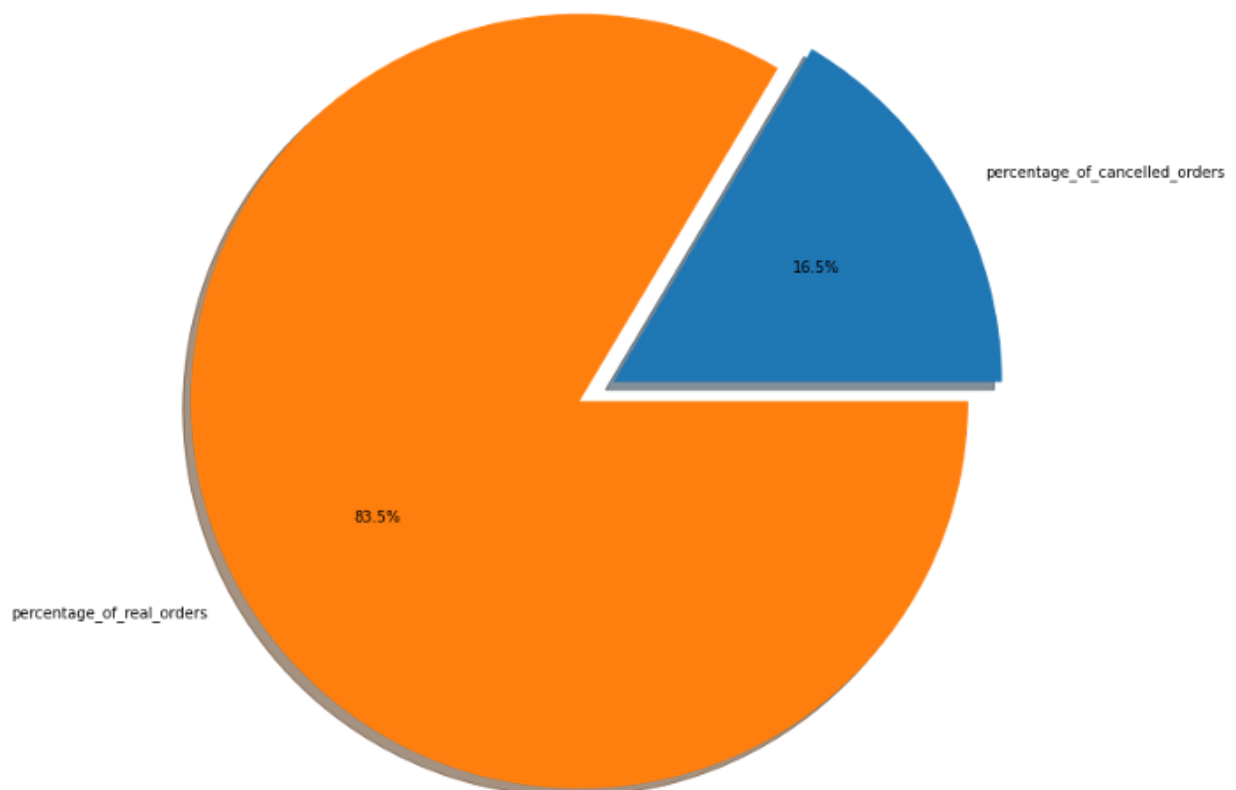
# Some Details we Fetched

## 1. What was the total revenue?

The total revenue was £ 8.73M

## 2. What is the percentage of canceled orders and real orders?



Real Orders Vs Cancelled Orders

## 3. What is the total revenue per month from December 2010 to December 2011?



Total Revenue per month from Dec 2010 to Dec 2011

| Month | Total Revenue |
| --- | --- |
| 2011-11 | 1,156,206 |
| 2010-12 | 570,423 |
| 2011-12 | 517,190 |
| 2011-10 | 1,035,642 |
| 2011-9 | 950,690 |
| 2011-5 | 677,355 |
| 2011-6 | 660,046 |
| 2011-8 | 644,051 |
| 2011-1 | 568,101 |
| 2011-4 | 468,374 |
| 2011-7 | 598,963 |
| 2011-3 | 594,082 |
| 2011-2 | 446,085 |

# 4. What are the top 10 products in terms of sales and revenue?

## Top 10 products in sales



SUM(Quantity)
26,076   80,995

- PAPER CRAFT , LITTLE BIRDIE - 80,995
- MEDIUM CERAMIC TOP STORAGE JAR - 77,916
- WORLD WAR 2 GLIDERS ASSTD DESIGNS - 54,319
- JUMBO BAG RED RETROSPOT - 46,078
- WHITE HANGING HEART T-LIGHT HOLDER - 36,706
- ASSORTED COLOUR BIRD ORNAMENT - 35,263
- PACK OF 72 RETROSPOT CAKE CASES - 33,670
- POPCORN HOLDER - 30,919
- RABBIT NIGHT LIGHT - 27,153
- MINI PAINT SET VINTAGE - 26,076

Quantity

## Top 10 products in revenue



SUM(Total Revenue)
42,584   168,470

- PAPER CRAFT , LITTLE BIRDIE 168,470
- REGENCY CAKESTAND 3 TIER 142,265
- WHITE HANGING HEART T-LIGHT HOLDER 100,392
- JUMBO BAG RED RETROSPOT 85,041
- MEDIUM CERAMIC TOP STORAGE JAR 81,417
- PARTY BUNTING 68,785
- ASSORTED COLOUR BIRD ORNAMENT 56,413
- RABBIT NIGHT LIGHT 51,251
- CHILLI LIGHTS 46,265
- PAPER CHAIN KIT 50'S CHRISTMAS 42,584

Total Revenue

# 5. Which products were returned more frequently?



**Top 10 cancelled products**

| Product | Quantity |
|---|---|
| PAPER CRAFT , LITTLE BIRDIE | 80,995 |
| MEDIUM CERAMIC TOP STORAGE JAR | 74,494 |
| ROTATING SILVER ANGELS T-LIGHT HLDR | 9,367 |
| FAIRY CAKE FLANNEL ASSORTED COLOUR | 3,150 |
| WHITE HANGING HEART T-LIGHT HOLDER | 2,578 |
| GIN + TONIC DIET METAL SIGN | 2,030 |
| HERB MARKER BASIL | 1,527 |
| FELTCRAFT DOLL MOLLY | 1,447 |
| TEA TIME PARTY BUNTING | 1,424 |
| PAPER POCKET TRAVELING FAN | 1,385 |

SUM(Quantity)
1,385          80,995

# 6. What are the top 10 countries that purchased the most?



**Top 10 countries by sales**

| Country | Quantity |
|---|---|
| United Kingdom | 4,247,040 |
| Netherlands | 200,834 |
| EIRE | 140,283 |
| Germany | 118,033 |
| France | 110,594 |
| Australia | 84,198 |
| Sweden | 36,037 |
| Switzerland | 29,981 |
| Spain | 27,735 |
| Japan | 26,016 |

# 7. Result after Clustering