



## A CLASSIFICATION APPROACH FOR NAÏVE BAYES OF ONLINE RETAILERS.

Aida Mustapha, Shazwani Mustapa, Nurfarahim Md.Azlan, Noor Fatin Ishmah Saifarrudin, Shahreen Kasim, Mohd Farhan Md Fudzee, Azizul Azhar Ramli, Hairulnizam Mahdin, Seah Choon Sen

Faculty of Computer Faculty Computer Science & Information Technology, University Tun Hussein Onn Malaysia, Johor, Malaysia  
seansea0702@gmail.com, {hana, shahreen, farhan, azizulr, hairuln}@uthm.edu.my

*This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.*

### ARTICLE DETAILS

#### Article history:

Received 22 January 2017  
Accepted 03 February 2017  
Available online 05 February 2017

#### Keywords:

consumer-centric marketing,  
online retailer

### ABSTRACT

Many small online retailers and new entrants to the online retail sector are keen to practice data mining and consumer-centric marketing in their businesses yet technically lack the necessary knowledge and expertise to do so. In this article a case study of using data mining techniques in customer-centric business intelligence for an online retailer is presented. The main purpose of this analysis is to help the business better understand its customers and therefore conduct customer-centric marketing more effectively. On the basis of the Recency, Frequency, and Monetary model, customers of the business have been segmented into various meaningful groups using the classification and naïve bayes algorithm, and the main characteristics of the consumers in each segment have been clearly identified. Accordingly a set of recommendations is further provided to the business on consumer-centric marketing.

### 1. Introduction

Data mining is the process of discovering interesting knowledge, such as associations, patterns, changes, significant structures and anomalies, from large amounts of data stored in databases or data warehouses or other information repositories. It has been widely used in recent years due to the availability of huge amounts of data in electronic form, and there is a need for turning such data into useful information and knowledge for large applications. These applications are found in fields such as Artificial Intelligence, Machine Learning, Market Analysis, Statistics and Database Systems, Business Management and Decision Support.

The online retailer under consideration in this dataset is a UK-based and registered non-store business with some 80 members of staff. The company was established in 1981 mainly selling unique all-occasion gifts. For years in the past, the merchant relied heavily on direct mailing catalogues, and orders were taken over phone calls. It was only 2 years ago that the company launched its own web site and shifted completely to the Web. Since then the company has maintained a steady and healthy number of customers from all parts of the United Kingdom and Europe, and has accumulated a huge amount of data about many customers. The company also uses Amazon.co.uk to market and sell its products [1].

The main purpose of this analysis is to help the business better understand its customers and therefore conduct customer-centric marketing more effectively. So, customers of the business have been segmented into various meaningful groups using the naïve bayes classification algorithm the main characteristics of the consumers in each segment have been clearly identified. Accordingly, a set of recommendations is provided to the business on customer-centric marketing and further data analysis tasks. The dataset was collected by using mark reports.

On the basis of the RFM model, customers of the business have been segmented into various meaningful groups using the clustering algorithm and decision tree induction, and the main characteristics of the consumers in each segment have been clearly identified.

### 2. Related Work

Researchers have proposed many different approaches for sentiment analysis. In this previous technical article for the online retailers industry is used k-means clustering to create some nested segments internally inside the cluster. In other words, these nested segments form some sub-clusters inside cluster, and make it possible to categorize the consumers concerned into some sensible subcategories.

With the prepared target dataset they intended to identify whether consumers can be segmented meaningfully in the view of recency, frequency and monetary values. The k-means clustering algorithm was employed for this purpose, and it can be easily performed by using the Cluster node in SAS Enterprise Miner. As well-known, the k-means clustering algorithm is very sensitive to a dataset that contains outliers (anomalies) or variables that are of incomparable scales or magnitudes.

### 3. Methodology

For this online retailer dataset, we choose classification methodology. Classification is a function of data mining, distributed collection in a project goal category or class. Classification is the purpose of accurately predicting each case of target class data. Classification task begins in a data set of class assignments are known Classification is a data mining function that assigns items in a collection to target categories or classes. The goal of classification is to accurately predict the target class for each case in the data. For example, a classification model could be used to identify loan applicants as low, medium, or high credit risks.

A classification task begins with a data set in which the class assignments are known. For example, a classification model that predicts credit risk could be developed based on observed data for many loan applicants over a period of time. In addition to the historical credit rating, the data might track employment history, home ownership or rental, years of residence, number and type of investments, and so on. Credit rating would be the target, the other attributes would be the predictors, and the data for each customer would constitute a case.

#### Advantages

- Users are thus better able to find precisely what they want, whether they wish to search in one discipline or across all. Works need not but can be coded by discipline.
- While other classification systems provide specific instructions in multiple places for coding by time or place or people, this system has a universal coding for such elements. This renders both classification and searching easier.
- The classification is able to use natural language but yet there are very precise definitions associated with each concept in the BCC. Following these steps a target dataset for the analysis has been generated. The original dataset was in MS Excel format, and was Part of the target dataset is shown in figure 1.

**Figure 1:** sample of the target dataset

```
graph TD; A[Training Data] --> B[Machine Learning Algorithms]; B --> C[Classifier]; D[New Data] --> C; C --> E[Prediction]
```

The flowchart illustrates the machine learning process. It begins with 'Training Data' (grey box) leading to 'Machine Learning Algorithms' (blue box). This leads to a 'Classifier' (grey box), which also receives input from 'New Data' (grey box). Finally, the 'Classifier' leads to the 'Prediction' (grey box).

**Figure 2:** A simplified diagram of the general model building procedure for pattern classification.

<b>Variable Name</b>	<b>Data Type</b>
Voice No	Nominal
Stock_Code	Numeric
Description	Nominal
Quantity	Numeric
Invoice_Date	Nominal
Unit_Price	Numeric
Customer_Id	Numeric
Country	Nominal

## 4. Experiment

The experiments that were carried out in this research were done using Rapidminer. Rapidminer is a java-based open source data mining and machine learning software. It has a graphical user interface (GUI) where the user can design his machine learning process without having to code. Here, the classification task for a given classifier is designed as a process. From the data we gathered, Algorithm that we use in online retailer is naïve bayes classification algorithm. We provided with a training dataset consisting of information about online retailers. Figure 3 show the dataset was in the form of a Microsoft Excel 2003 spreadsheet and had details of country, invoice number, stock code, quantity, description, quantity, invoice date, unit price and customer id. For ease of performing data mining operations, the data was change into a csv format.

Figure 3: Dataset for Online Retailers

### b. Algorithm of naive Bayesian Classifier

The screenshot displays the Proteomix software interface. The main workspace shows a workflow with three blocks: 'Select Regions' (green), 'Apply Model' (blue), and 'Performance' (yellow). The 'Performance' block is selected, and its parameters are visible on the right. The parameters include: 'Validation (X Validation)' checked, 'Average performance only' selected, 'Leave one out' unselected, 'Number of validations' set to 10, and 'Sampling type' set to 'Automatic'.

**Figure 4 : Design view of Naive Bayes**

Performance (Per Performance)					ExampleSet (Set Role)		SimpleDistribution (Name Bayes)	
ExampleSet (999 examples, 1 special attribute, 7 regular attributes)					Filter (999 / 999 examples)			
Item No.	Country	InvoiceNo	StockCode	Description	Quantity	InvoiceDate	UnitPrice	CustomerID
1	United Kingdom	536365	851234	WHITE HANGG.	6	01/12/2010	2.300	17890
2	United Kingdom	536365	71053	WHITE METAL	6	01/12/2010	2.350	17890
3	United Kingdom	536365	644008	KREAM CUPR.	6	01/12/2010	2.750	17890
4	United Kingdom	536365	840290	ONION LENT.	6	01/12/2010	3.300	17890
5	United Kingdom	536365	840936	RED WIGGL.	6	01/12/2010	5.100	17890
6	United Kingdom	536365	29709	SET 2 BULL.	6	01/12/2010	7.650	17890
7	United Kingdom	536365	21730	GLASS STAB.	6	01/12/2010	9.100	17890
8	United Kingdom	536365	25033	HAND WASHB.	6	01/12/2010	1.850	17890
9	United Kingdom	536365	22632	HAND WASHB.	6	01/12/2010	1.850	17890
10	United Kingdom	536367	84479	ASSORTED	52	01/12/2010	1.600	13047
11	United Kingdom	536367	22745	POPPYS PLA.	6	01/12/2010	2.100	13047
12	United Kingdom	536367	22748	POPPYS PLA.	6	01/12/2010	2.100	13047
13	United Kingdom	536367	22749	FELT CRAFT	6	01/12/2010	3.750	13047
14	United Kingdom	536367	22310	HORT KENIT.	6	01/12/2010	1.650	13047
15	United Kingdom	536367	84899	BOX OF 8 AG.	5	01/12/2010	4.250	13047
16	United Kingdom	536367	22623	BOX OF VINT.	5	01/12/2010	4.900	13047
17	United Kingdom	536367	22502	BOX OF WINT.	2	01/12/2010	9.950	13047

**Figure 5 : Set role of Online Retailers**

Evaluation Metric is a process of assessing progress towards the goals and objectives that have been defined. It including information on the efficient use of resources in producing to compare with the results of the intended meaning and effectiveness in achieving the objectives of the action. This evaluation metric will be proof to show that the online retailer classification are correctly or incorrectly classified using classification accuracy, coefficient metric and time taken test model on dataset. The Confusion Matrix obtained is given below in instances and 8 attributes is fed as input to the Naïve bayes figure 6 and figure 7 is plot view for analysis.

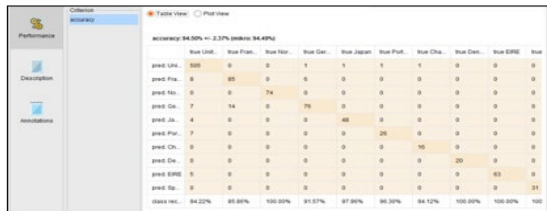


Figure 6 : Naïve Bayes Confusion matrix

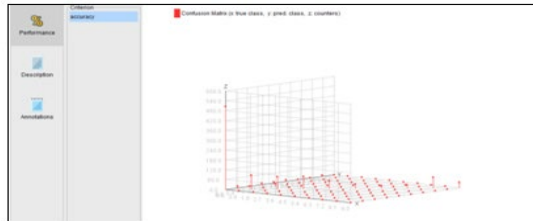


Figure 7: Plot View

## 5. Result

Results of the empirical study show that prediction the country in online retailer dataset. The classification of this datasets was conducted to prove that the histograms in figure 8 of the variables country the target dataset of online retailer it is evident that there are a few instances having quite different monetary and frequency values compared to the majority of the country in the dataset. These instances are valid from the business point of view as they are genuine transaction records; however, they are outliers from the data analysis point of view. United kingdom is the highest ranking country for delivery the item.

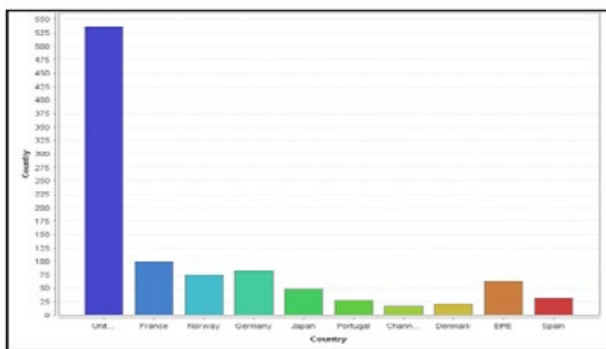


Figure 8: Bar chart of country

RapidMiner helped significantly in finding hidden information from the training dataset. These newly learnt predictive patterns for predicting country performance in online retail. In figure 9 it show the simple distribution for label attribute country and figure 10 show the distribution table. We can see each country have 7 distribution.



Figure 9: Simple Distribution for Country

Simple distribution included distribution model for label attribute Country which is United Kingdom, France, Norway, Germany, Japan, Portugal, Chanel Island and Denmark.

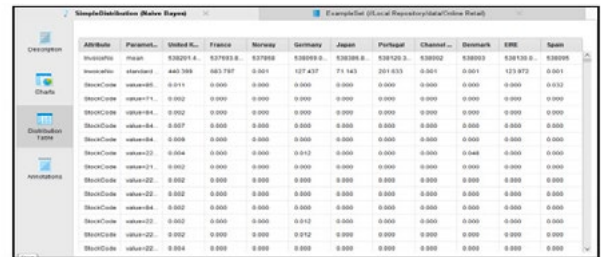


Figure 10: Distribution table

In figure 11, it shows the example chart of country based on the quantities of each product (item) per transaction. We can see that portugal is the highest density with 0.068.

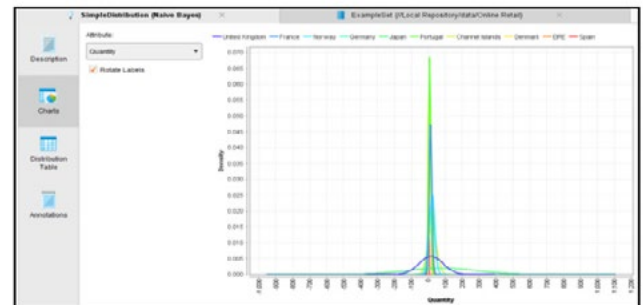


Figure 11 : Chart of Country based on Quantity

## 6. Conclusion

A case study has been presented in this article to demonstrate how customer-centric business intelligence for online retailers can be created by naïve bayes of data mining techniques. The distinct customer groups characterized in the case study can help the business better understand its customers in terms of their profitability, and accordingly, adopt appropriate marketing strategies for different consumers. It has been shown in this analysis that there are two steps in the whole datamining process that are very crucial and the most time-consuming: data preparation and model interpretation and evaluation.

## 7. References

1. Kumar , V . and Reinartz , W . J . ( 2006 ) Customer Relationship Management: A Databased Approach ,Hoboken, NJ: John Wiley & Sons .
2. "A tutorial on support vector machines for pattern recognition," Data Mining and Knowledge Discovery, vol. 2, no. 2, pp. 121-167, 1998 by C. Burges.
3. Bhardwaj, B. K., & Pal, S. (2012). Data Mining: A prediction for performance improvement using classification. arXiv preprint arXiv:1201.3418.
4. Sebastian Raschka (2013) Naive Bayes and Text Classification