

An in-depth investigation across different product sales

Keith Conti Borda

Institute of Information & Communication Technology

University College

MCAST, Paola PLA 9032

keith.conti.a100360@mcast.edu.mt

Abstract—Investigating across the different sales of various products in different areas by different individuals. The aim of this project is to analyze the different buying behaviours, to try and identify certain buying patterns for different products, buyers and also supply stores from the data set which has been provided to us. The data set consists of a combination of categorized products, orders, stores and clients, from which we have to extract the data and insert into our customized cube. The proposed solution includes a custom OLAP cube which has been implemented to cater for the data stored in the original data set after being cleansed for duplicates. This project's outcome has resulted in a positive one, since all the migrations of data have been executed successfully. A Windows based system was used for this project evaluation, together with SQL Server Management Studio, Microsoft Excel and the Oracle Data Modeller to design the ERD's.

Index Terms—analysis, investigation, dataset, data cleansing, operations

I. INTRODUCTION

With regards to this whole project, alot of effort and energy has been put in by our tutor, in order to provide us with the dataset, and all the necessary material in in order for us to be able to work and implement the required operations for this scenario. The aim of this entire project was to analyze the different buying behaviours, to try and identify certain buying patterns for different products, buyers and also supply stores from the data set which has been provided to us by our tutor. Having this data stored in the dataset, we were required to compile our own SQL Scripts in order to query the dataset, to collect data and information we need in order to satisfy the requested reports. By having done this, we were able to identify who bought which product, from which category, from which store in which region, during what year, month, quarter and even day. All the data provided to use by our tutor inside the olap database, was filtered thoroughly, and then was migrated into our own custom OLAP cube, whereby the data may be displayed in a similar manner, however cleansed.

II. LITERATURE REVIEW

It has become a vital and extremely challenging for organizations to sustain their position within a competitive market and also to understand the customer. Most importantly, technological improvements have paved a way in order to process queries rapidly and efficiently, thereby minimizing the waiting time. Data mining utilities have become the most

important tools in order to analyze and filter the huge amounts of data to make the correct decisions. The objective of this paper is to analyze big data together with the consumer's behaviours and also to make the correct decisions leading to competitive rivals.

The database technology we know of today has been changed drastically since the 1980s. Research and development activities on new and powerful database systems has advanced much more than it used to be, thereby employing the advanced data model. The extreme growth of computer systems, together with the hardware and software technology in the past three decades has led to a large supply of powerful, and and cost effective computer systems, data collection equipment and also storage media. Thanks to all this, a huge number of databases and data repositories were made available for transaction management information retrieval and data analysis. The online analytical processing (OLAP) withholds analysis techniques with functionalities such as summarizing, consolidation and aggregation together with the view of data from different angles. These OLAP tools have been used for in depth analysis such as classifications of data, clustering and the data changes over time. The term 'Data mining' refers to the extract of knowledge from huge amounts of data. This is a simple step in the knowledge process to extract knowledge store inside databases.

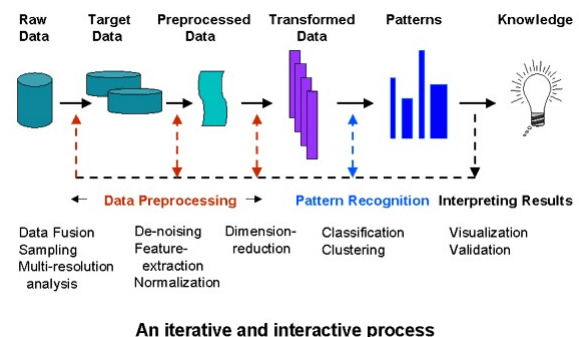


Fig. 1. DM Diagram

The forecasting and sales analyses of certain products and services makes a difference in the profit and loss per accounting period for the organization which then determines

the success or failure of the organization itself. Therefore, the reliable prediction of sales is of utmost importance. This article presents the sales forecasting approach by the integration of fuzzy systems and data clustering to construct a sales forecasting integrated system. All clusters are entered into an independent GFS model with the ability of rule base extraction and database tuning. The performance of an extracted expert system is compared against the previous sales forecasting methods using mean absolute percentage error (MAPE) and root mean square error (RMSE). These results show that the proposed approach outperforms the other previous approaches.

Market sales is the only factor than measures the outcome of marketing efforts and therefore the sales forecasting are extremely important when it comes to marketing planning. Within the competitive market each business faces, sales forecasting is a very helpful tool, because an organization may use this to its advantage to have a rough idea of what sales to expect based on current of past figures. A reliable sales predication will improve the organizations sales because it entitles production, sales, marketing, and finance sections to effectively develop programs to continue enhancing the organizations size. Examples of these software programs are: Sales planning, budgeting, promotion and advertising plans.

Market basket analysis is also known as association-rule mining. This terminology is a very useful method of discovering customer purchasing patterns data from stores transactional databases. Since the information obtained from the analysis may be utilized to form marketing, sales, service, and operation strategies, it increased the research interest. However, the existing methods may fail to discover important purchasing patterns in a multi-store environment, because of an assumption that products under consideration are on shelf all the time across all stores. Within this paper, a new method is being proposed to overcome this weakness. This evaluation shows that the proposed method is efficient, and has an advantage over the traditional method when stores are varied in size, product mix changes rapidly over time, and larger numbers of stores and periods are considered.

Thanks to all the advances in information technology we have today, organizations may effectively gather and store transactional and demographic data on individual customers at very reasonable costs. The most popular challenge for all corporations is that of how to extract customer data information from their vast customer databases to gain a competitive advantage. This Market Basket analysis is a method of how to discover customer purchasing patterns by extracting data from their transactional databases. The result of this method may be rather biased since a product may be on shelf before its first transaction or also after the last transaction. This product may also be put on shelf and taken off shelf many times during this market basket data collection period. This is the first problem within this research. The second problem is that

of finding the common association patterns in subsets of stores. To overcome these problems, an algorithm was developed in order to automatically extract association rules in a multi store environment, however the rules also contain the information on store together with the time where the rules hold.

Since the time and store location factors will be considered, the rule generation is much more complicated than the algorithm mentioned above. Within this paper, the simulation results which show the proposed method that is computationally efficient and has an advantage over the traditional association method, when stores under observation are large in size and have large mixes of products that are altered rapidly over a period of time.

III. RESEARCH METHODOLOGY

The CRISP DM terminology stands for cross-industry process for data mining. This methodology offers an approach that is structured to plan a data mining project. The crisp-dm methodology is powerful, practical, flexible and useful when utilized to its advantage to solve business issues. The CRISP-DM model is shown hereunder. [1]

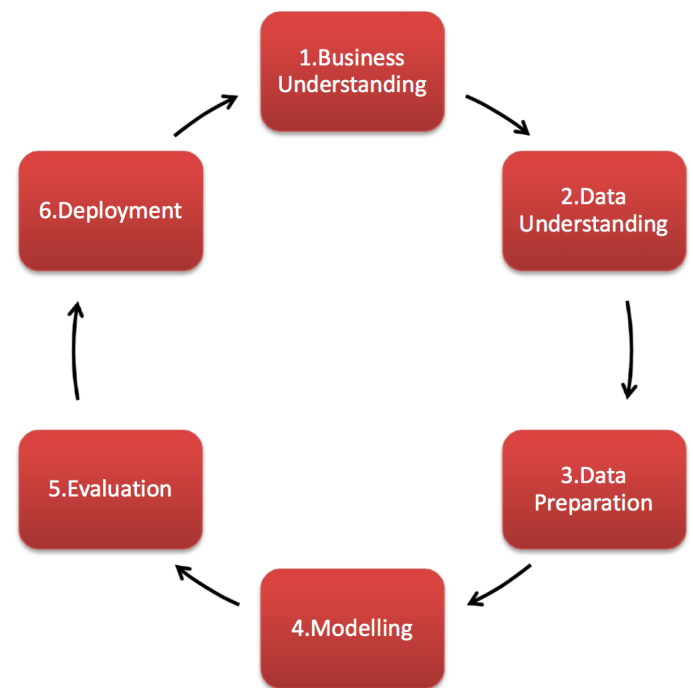


Fig. 2. crispDM

This model is made up of a sequence of stages, which may be performed in different order and will often be required to backtrack to previous tasks and repeat certain actions. This model does not capture all the routes through the data mining process. The 6 stages that make up the crisp-dm are as follows: Business Understanding
Within the first stage of the crisp-dm, one has to understand what has to be accomplished from a business perspective,

and after this happens the knowledge is converted into a data mining problem definition. The goal of this stage is to uncover important factors which may influence the outcome of the project. Ignoring this step may drop all the effort made thereby producing the right answers to the wrong questions. Data Understanding

The second stage of the crisp-dm starts with an initial data collection and proceeds with activities in order to familiarize along with the data, to identify the problems with data quality, find first insights in the data, or even to detect interesting subsets to form hypothesis for hidden information. Data Preparation

The third stage of the crisp-dm is the data preparation. This stage covers all activities to build the finalized dataset. These preparation tasks are performed multiple times, and not in any prescribed order. These tasks mentioned include tables, records and also attribute sections together with transformation and the cleaning of data for the modelling tools. Modelling

The fourth stage of the crisp-dm is that of modelling. There are various modelling techniques that are chosen. Typically, there are many techniques for the same data mining problem type. Some of the techniques have specific requirements when it comes to the form of the data, thereby we would have to step back to the data preparation phase. Evaluation

Within the fifth stage of the crisp-dm methodology, the model built will appear to have a very high quality, from a data analysis perspective. Before advancing to the final deployment of the model, it is very important to have an even deeper evaluation on the model itself, and after that review the steps used to build the model to be sure it achieves the business objectives specified. At the end of this phase, a decision on the use of the data mining results must be reached. [2]

Deployment

Finally, the sixth stage of the crisp-dm methodology. The creation of the model is not the end of the project. Even though the purpose of the model improves the knowledge of the data, the knowledge gained must be organized and shown in a certain way that would be beneficial to the customer. Depending on the requirements, the deployment phase can be as simple as generating a report or as complex as implementing a data mining process. In many cases, it will be the customer, not the analyst, carries out the deployment steps. Even if the analyst deploys the model it is important for the customer to have a good understanding of the actions which will need to be carried out in order to actually make use of these created models.

In order to be able to design the 3D OLAP cube, we had to first create our own ERD from the original OLTP ERD. In my case, my erd consisted of 5 tables including the Order Fact table. Once this was complete, I then designed the 3D OLAP as per below diagram.

The OLTP ERD provided to us by our tutor shows all the tables that may be found inside the database, including all their rows and data types, primary keys, foreign keys in addition with the way they are linked to each other. This gives us a

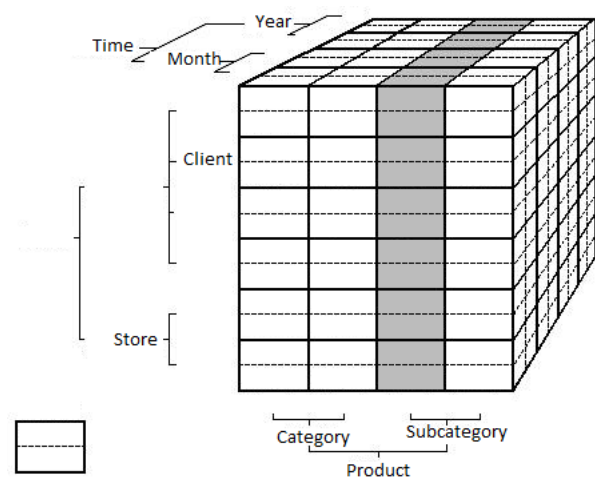


Fig. 3. OLAPTemplate

better understanding of how we can aggregate the data from more than one table to another. The original OLTP ERD may be seen hereunder;

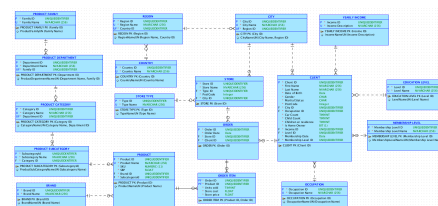


Fig. 4. Original OLTP ERD

What we had to do was the following; We created our own ERD, by choosing a dimension of our choice. In my case, my fact table was the Orders fact, which is connected to another 4 four tables being Client, Store, Time and Product. Obviously, the tables I chose were taken from the original ERD (OLAP), together with their respective columns and data types too. This entitles us to query and display the data from another dimension; thereby we can say what type of client bought which product, from which store in which region, during what year, month or even quarter and also on which day. The ERD I have created may be found hereunder;

With regards to the preliminary analysis and data cleaning, a script was built in order to gather all those duplicate records inside the OrderItem table inside the database.

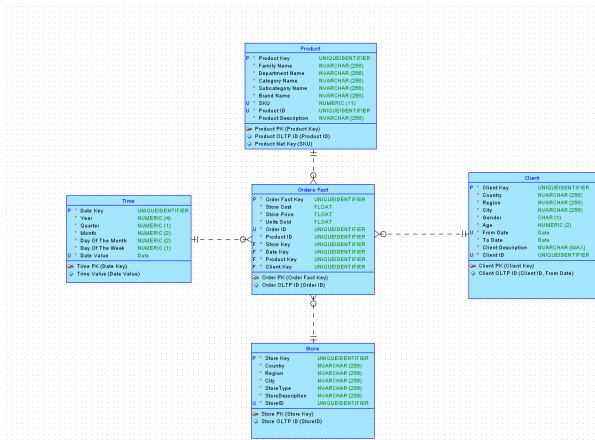


Fig. 5. Order Fact ERD Dimension

```

10 -- Preliminary Analysis to check for duplicates
11 SELECT
12     productId,
13     orderId,
14     COUNT(*) AS 'Count'
15 FROM
16     [oltp].[orderItem]
17 GROUP BY
18     productId, orderId
19 HAVING COUNT (*) > 1

```

Fig. 6. Code For Duplicates

The duplicates found by executing the above SQL query were as follows:

Once these duplicates were analyzed, another SQL Query was built in order to view only those duplicate records inside the OrderItem table in the database.

By making use of a Common Table Expression (CTE), the duplicate records inside the database were removed, thereby only having unique values in the table without any duplicated values.

ETL is used to migrate data from one database to another, to form data marts and data warehouses and also to convert databases from one format or type to another. ETL is an acronym which stands for Extract, Transform and load. These three are database functions that are combined into one tool and are used to extract information from the database, and insert it into another database or table.

Extract the process of reading the data stores inside a table in the database.

Transform the process of converting the gathered data from extraction into the form it needs to be stored in so that it may be placed inside another table or database. The transformation happens by using the lookup tables or by combining the data with the other data.

Load the final process of the ETL. The data is written into the target table or database.

In our scenario, we made use of the ETL procedure in order to migrate the stored data given to us originally into the OLAP 3D cube we designed ourselves.

After the ETL implementation and execution was complete, the data was migrated from the original OLTP and inserted into the OLAP cube created by us. In order to confirm that

BI Assignment Scri...RDA-PC\keith (52)* x SQLQuery1.sql - KC...RDA-PC\keith (53))

```

10 -- Preliminary Analysis
11 SELECT
12     productId,
13     orderId,
14     COUNT(*) AS 'Count'
15 FROM
16     [oltp].[orderItem]
17 GROUP BY
18     productId, orderId

```

productId	orderId	Count
1	2	2
2	2	2
3	2	2
4	2	2
5	2	2
6	2	2
7	2	2
8	2	2
9	2	2
10	2	2
11	2	2
12	2	2
13	2	2
14	2	2
15	2	2
16	2	2
17	2	2
18	2	2
19	2	2
20	2	2
21	2	2
22	2	2
23	2	2
24	2	2
25	2	2
26	2	2
27	2	2
28	2	2
29	2	2
30	2	2
31	2	2
32	2	2
33	2	2
34	2	2
35	2	2
36	2	2

Fig. 7. Duplicates

```

18 -- CTE To Delete the duplicates found in the previous statement
19 WITH CTE
20 AS (SELECT productId, orderId, ROW_NUMBER() OVER (PARTITION BY orderId ORDER BY productId)
21 FROM [oltp].[orderItem])
22 DELETE FROM CTE
23 WHERE rowNumber > 1
24
25 -- Drop Tables
26 DROP TABLE [ProductFact].[Product];
27 DROP TABLE [ProductFact].[Client];
28 DROP TABLE [ProductFact].[Store];
29 DROP TABLE [ProductFact].[Time];
30 DROP TABLE [ProductFact].[OrderFact];
31 GO

```

Messages

(196795 row(s) affected)

Fig. 8. Common Table Expression

the data was stored successfully inside the new OLAP table within the database, we executed simple select statements to display the data that had just been stored inside the database.

IV. CONCLUSION

Thanks to this interesting project, us students / researchers have learnt alot about how the stored data may be extracted, analysed, transformed, migrated and stored inside a custom OLAP 3D cube by making use of ETLs. We were well equipped with study material in order to be able to work on



Fig. 9. ETL

this project and this may motivate us to keep on extending this study further.



Fig. 10. MCAST Logo

REFERENCES

- [1] “What is the crisp-dm methodology,” 2017. [Online]. Available: <http://www.sv-europe.com/crisp-dm-methodology/dataunderstanding>
- [2] “Cross industry standard process for data mining,” 2017. [Online]. Available: <https://en.wikipedia.org/wiki/CrossIndustryStandardProcessforDataMining>