# Data Mining Techniques For the Integrated Postsecondary Data System

Wilfredo Jiménez Vargas

Jeffrey Duffany, Ph.D.

Electrical & Computer Engineering/Computer Science Department

## Abstract

The Process of generating the necessary information for the Integrated Postsecondary Data System can be tedious, a standard way to collect the information for the different department's needs.

To process it, we use Data Mining techniques to obtain consistent and replicable results that can be consumed by the end user for multiple criteria analysis that often reveal patterns that shape the future organizational choices. Some of the different techniques and concepts used will be discussed in further detail on this paper.

## Introduction

The Integrated Postsecondary Education System is a collection of data from higher education institutions that are accredited by the United States Department of Education. The completion of this report is mandatory of all institutions that take part in any federal assistance program authorized the title IV of the Higher Education Act of 1965. Usually, the data requested for this kind of report involves a population 'Cohort.' This population will be used to verify data; one example is the average amount of grants or the average amount of time a student finished their degree in the first 150 percent of the time.

## Background

Usually, the steps needed to collect and analyze the data required for this kind of report is a challenging and intricate process. This information is verified with previous year submissions and cross-checked with the current year data; another challenge is that the majority of this information is collected manually or requested to the Information System office in multiple reports generated by different actors.

## Problem

This kind of data collection fragmentation creates problems when the data is loaded to the system, and some of it needed to be explained. One central data collection for the distinct institutional resources is required to preserve the integrity and the traceability of the information.

## Methodology

In terms of data processing the traditional statistical models, have restrictions: This kind of processes require assumptions sufficient familiarity of probabilities and distributions, however, the information must have the prerequisite of being of high quality, being a target to prior processing and transformation. The very nature of this process poses a disadvantage, giving space to OLAP, which stands for (Online Analytical Processing) this, is a method by which multidimensional analysis works.

The data must be sourced from the different institutional departments, this implies that a process for data requirements gathering that must include various meetings with the department directors or delegates where arranged in order to understand the relationship between the different sources of information that will be used in order to create an accurate schema design. The relationship that the information have with different departments must be fully analyzed in order to design robust schemas and produce the necessary reports, this must be done in order to certify the information is accurate.

### What is MOLAP

Multidimensional OLAP (MOLAP) is a implementation of OLAP that facilitates data analysis by using a multidimensional data cube. Data is precomputed, pre-summarized, and stored. Using this method a stakeholder can use the multidimensional view to observe the data with different facets. The multidimensional data analysis can also be implemented in a relational database by querying multiple tables however this the version that is used for this project is the MOLAP because it has all the possible combinations of data already stored in a multidimensional array.
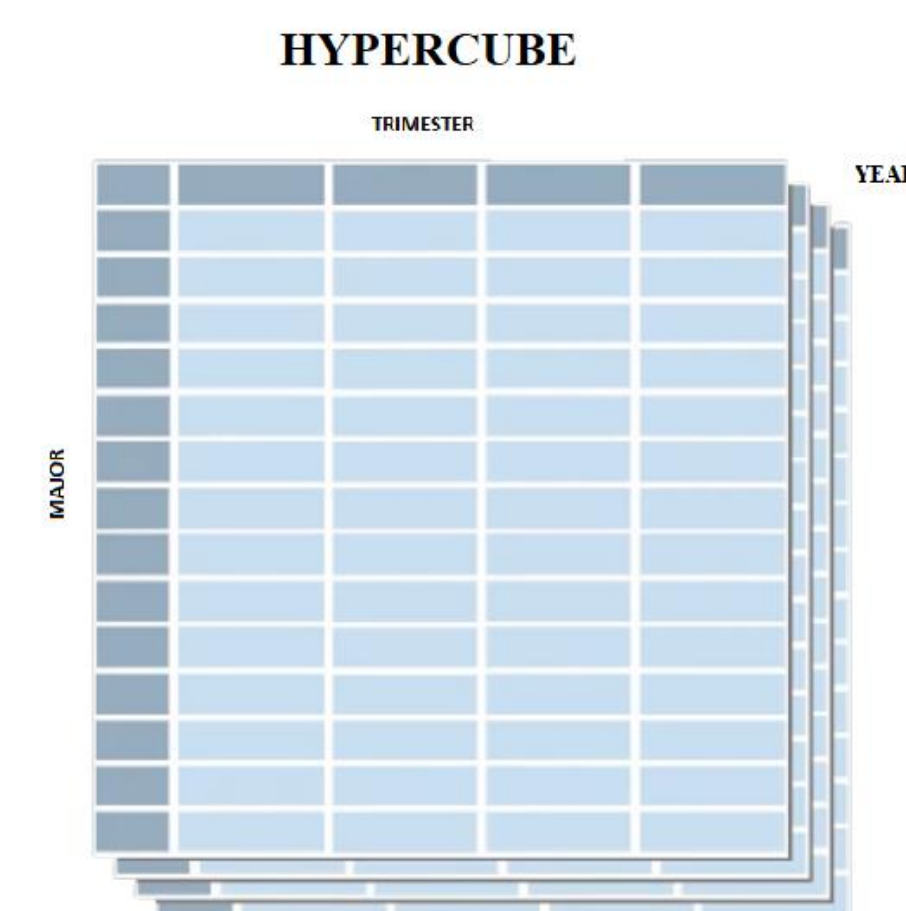


**Figure 1**
**IPEDS Hypercube**

### Key Points

- Using MOLAP to process information with the same time irrespective of the level of summarizing.
- It also remove the need for design a relational database to store data for analysis.
- Facts are stored in multidimensional arrays and dimensions used to query them.

## Results and Discussion

The complex consequences of certain decision-making situations necessitate taking into account the multicriteria and conflicting aspects of data, as well as the consideration of several types of data (quantitative and qualitative) to represent all the necessary information for making decisions adequately.



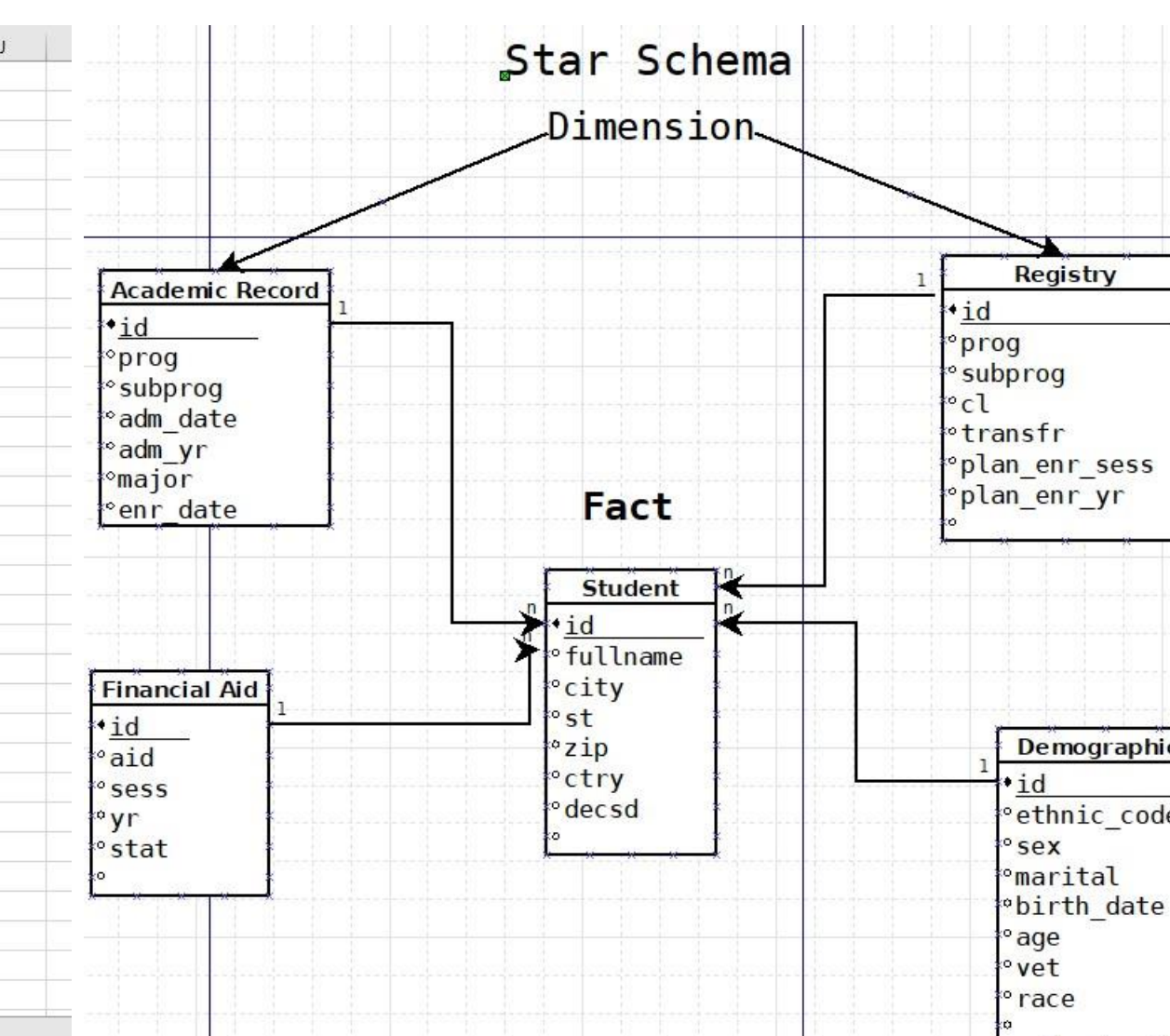**Figure 2**
**Data for Multidimensional Analysis**



**Figure 3**
**Star Schema**

However, OLAP tools that exist in the decision-making area still suffer from limitations related to the lack of technical means to consider the multiple criteria and imprecise nature of decision data in the analysis process.

In fact, OLAP systems are the cornerstone of many analysis applications and present as a complete system providing useful and necessary services for efficient, rational, and analytic processing of data.



**Figure 4**
**Modeled And Filtered Data**

## Conclusions

In overall, the success of this kind of data collection techniques include exhaustive communication with the stakeholder and the data presentation must be in a manner that the end user can relate. As a result, the Microsoft Excel was selected to display the information, because the familiarity and wide spread use prevented that end user felt alienated.

It is important to have in mind the necessities that the user might have, because they are the subject manner expert in their field and there is a need that the data is understandable and useful by them in order to be verified and certified by the departments.

## Future Work

At this point is uncertain what will be the changes for the next year reporting period. In addition, that sometimes the data source might change because of a modification on the system requirements, if this is the case the schemas must be updated and documented accordingly preventing future problems

## Acknowledgements

## References

[1]Codd, Edgar F., et al., "Beyond decision support," in Computerworld. Oxford University Press, [online document], 1993. Available: Ebsco Host Online, http://ezproxy.pupr.edu:2055/login.aspx?direct=true&db=bth&AN=9311083206&site=ehost-live [Accessed: Jan 24, 2019].

[2]Ralph Kimball, *The Data Warehouse Toolkit: Practical Techniques for Building Dimensional Data Warehouses*, 3rd ed. Reading, MA: Addison Wesley, 2013. [E-book] Available: Kindle e-book.

[3]BM Community," Dic, 2015. [Online]. Available: https://www.ibm.com/developerworks/community/blogs/2fa81a5c-cb30-4873-b775-1370151e3614?lang=en.htm. [Accessed Jan. 20, 2019].

[4]National Center for Education Statistics ," Jan, 2017. [Online].Available:https://nces.ed.gov/collegenavigator/?q=politecnica&s=all&id=243577. [Accessed Jan. 22, 2019]