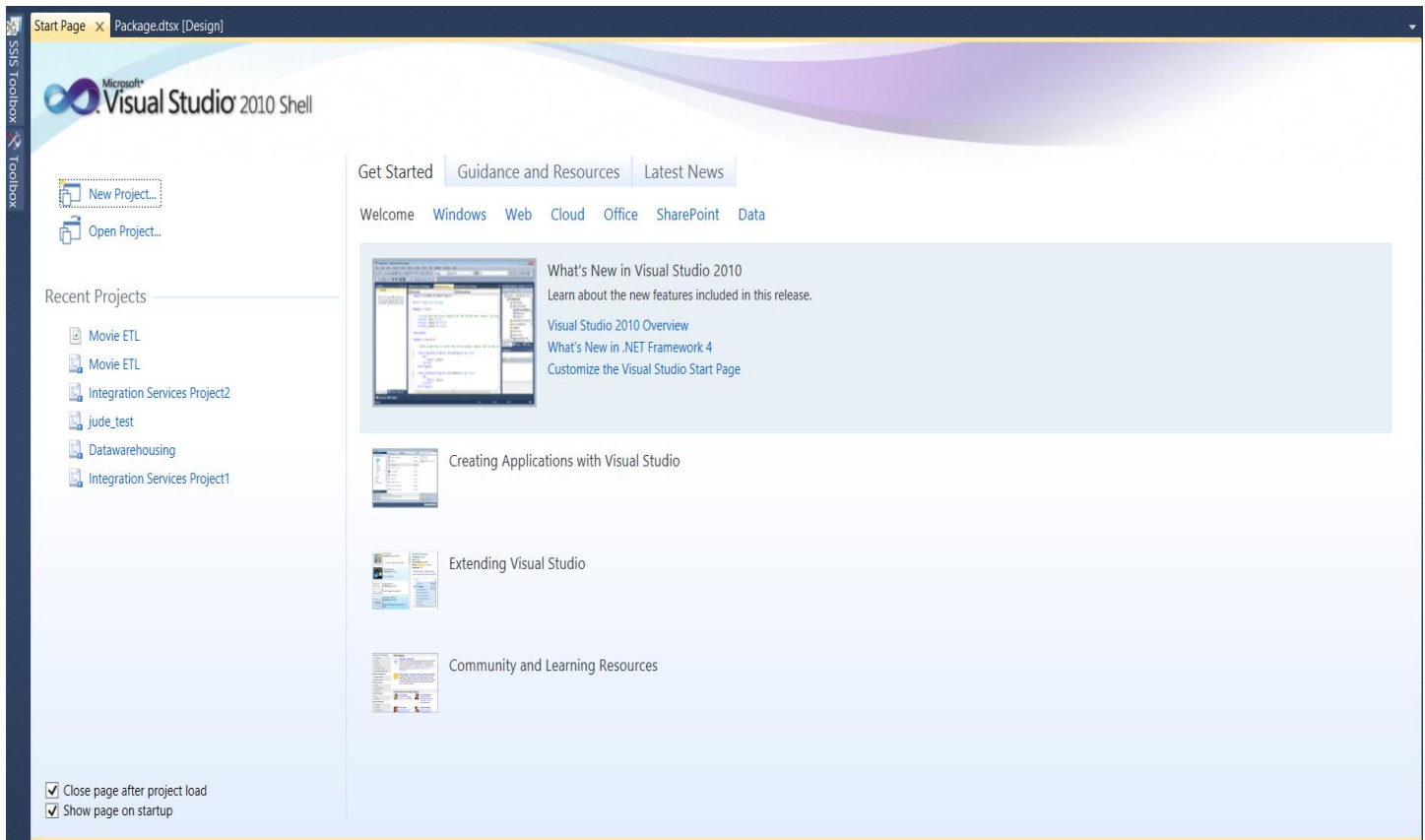# "Movie Management System"
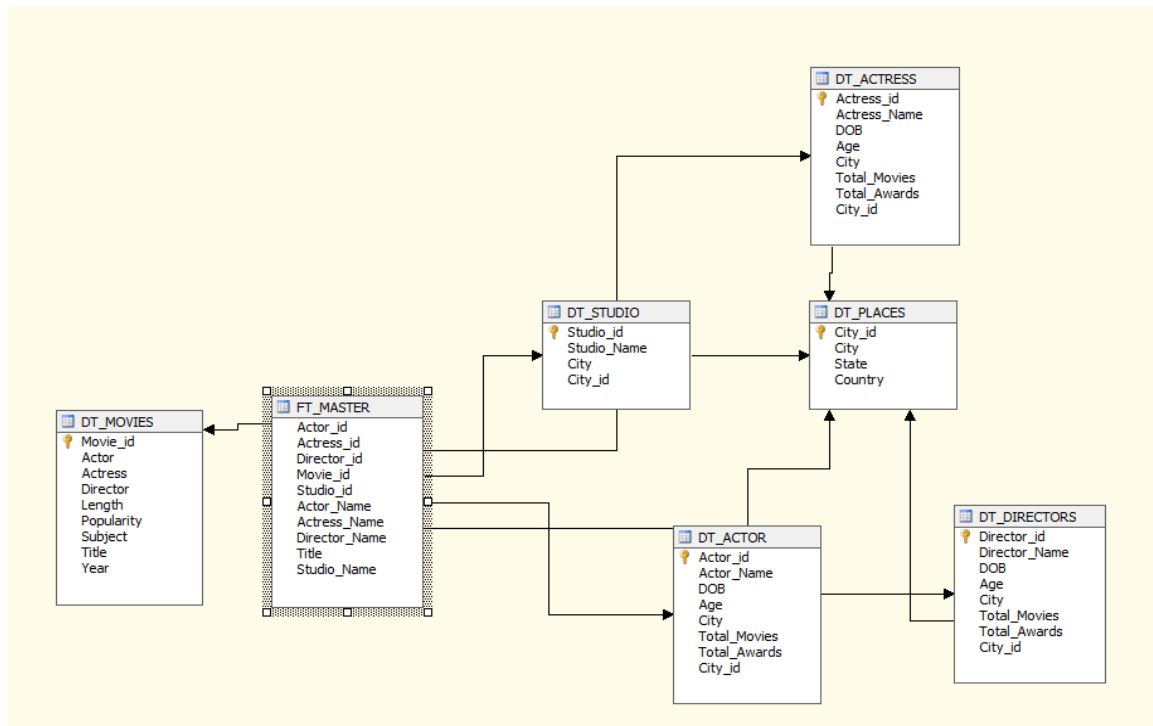## Data Mining of the Application

*Mohamed  Niyaz*

# 1.  Mining Introduction

Data mining consists of a set of statistical techniques for analyzing observational data. The techniques are employed in the analysis step of the "Knowledge Discovery in Databases" process, or KDD. These tools discover patterns in large data sets. The overall goal of the data mining process is to extract information from a data set and transform it into an understandable structure for further use. Aside from the raw analysis step, it involves database and data management aspects, data preprocessing, model and inference considerations, interestingness metrics, complexity considerations, post-processing of discovered structures, visualization, and online updating.

# 2.  Mining Design



Our Data Mining Process is carried through **Microsoft BI** tool and below screen shot explains the high level overview of the **Schema**

The below Screenshot and process depicts different parts of Data Mining Process.

**Tools Used: Microsoft SSAS**

## 3. A Description of our mining scenarios and implementation

We designed and implemented two mining scenarios for this phase:

**Scenario 1:**
The movie popularity system used by the IMDB uses a simple 0-100 scale. These values are given definite values for example 5, 10 where 0 being the minimum and 100 being the maximum. Therefore, in order to better understand this popularity system we want to cluster the movies by popularities. This will allow us to find breaking points between clusters and use these to apply labels to certain number ranges. In other words we hope to use clustering to identify clusters of movie ratings and assign values to each cluster.
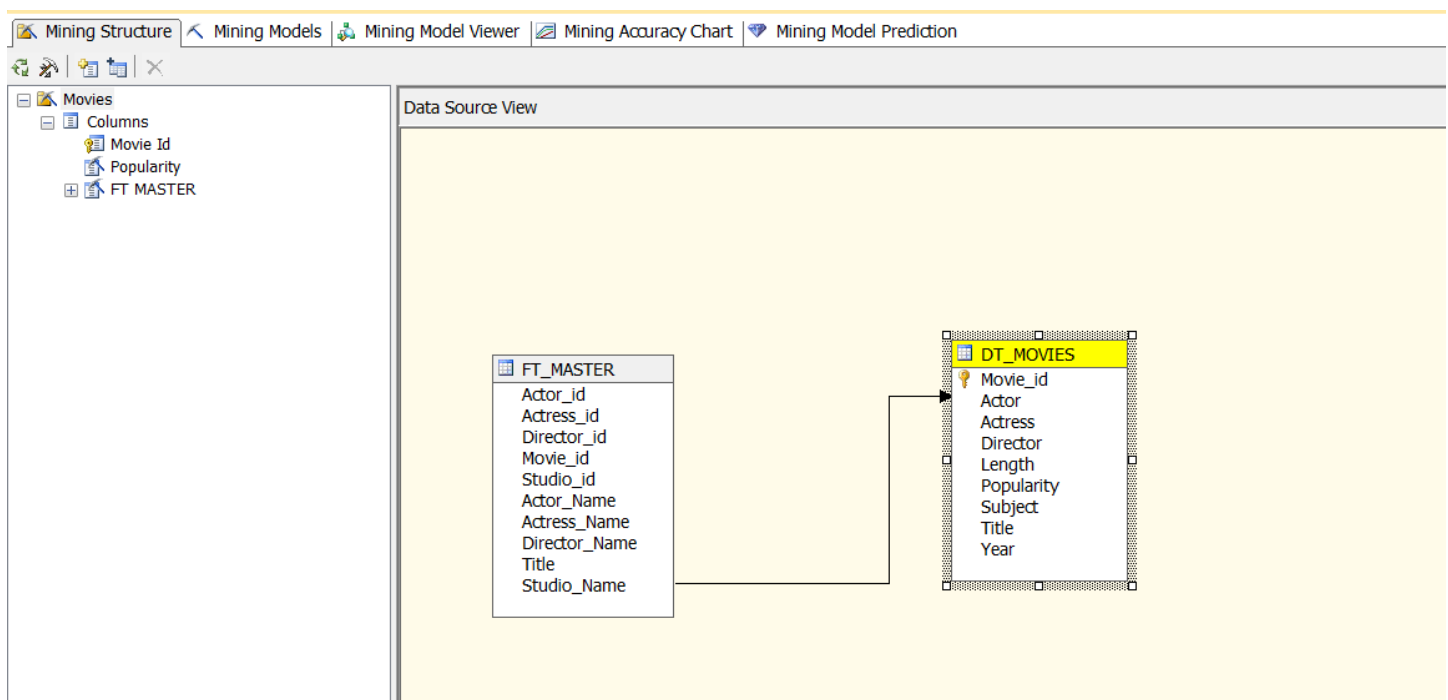
## 4. Implementation

We implemented mining scenario using Sql server analysis service.
Steps Involved:
1. Created a new SSAS project.
2. Gave connection to Sql Server 2012 and added an existing Movies database to the data source.
3. Created Data source view from the dimension and fact tables present in the database.
4. Created a cube structure using the fact table.
5. Created Mining structure, we used "Microsoft clustering algorithm". Added "DT_MOVIES" table as "Case" table and "FT_MASTER" table as the "nested" one.
6. Gave "rating" as prediction attribute and "Movie_id" as input attribute and "Actor_Id"as the key.

### *Mining Structure for the clustering Algorithm*

Overall relation of a Case table and a Nested Table

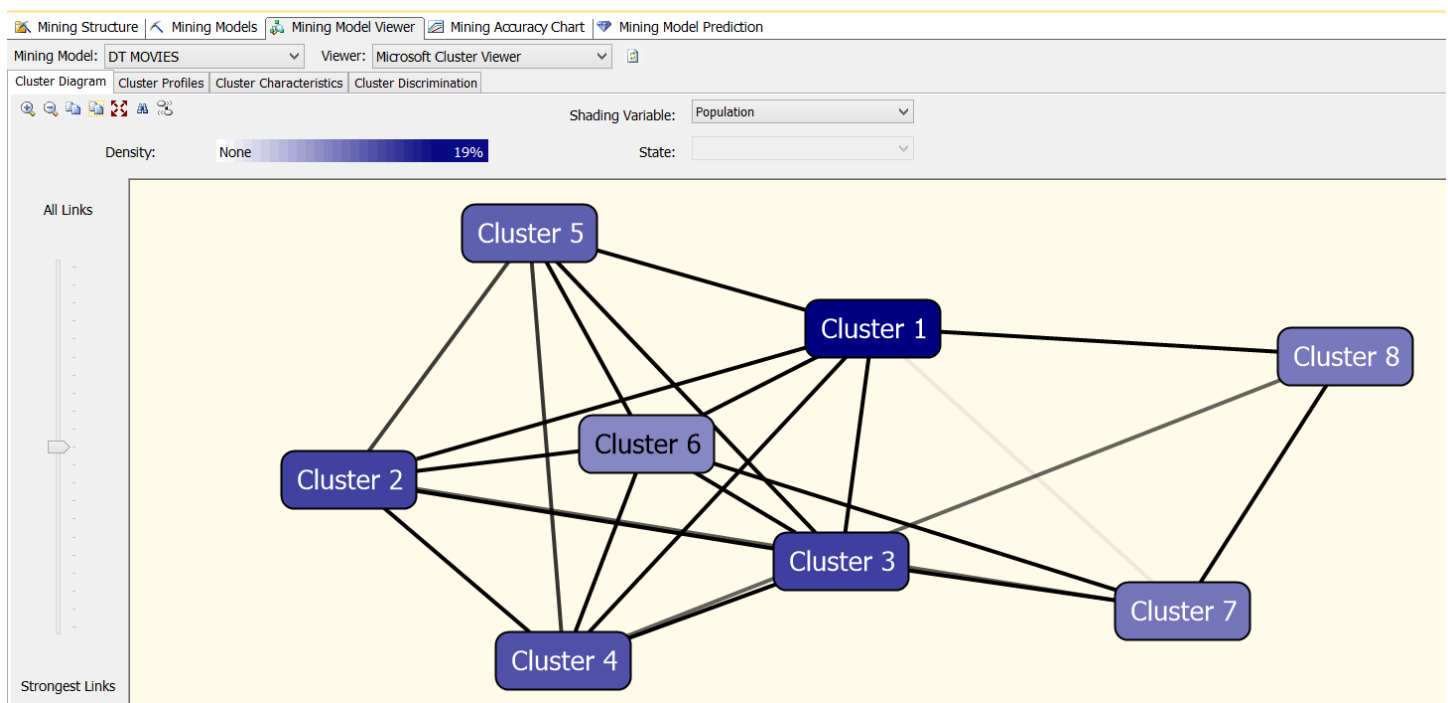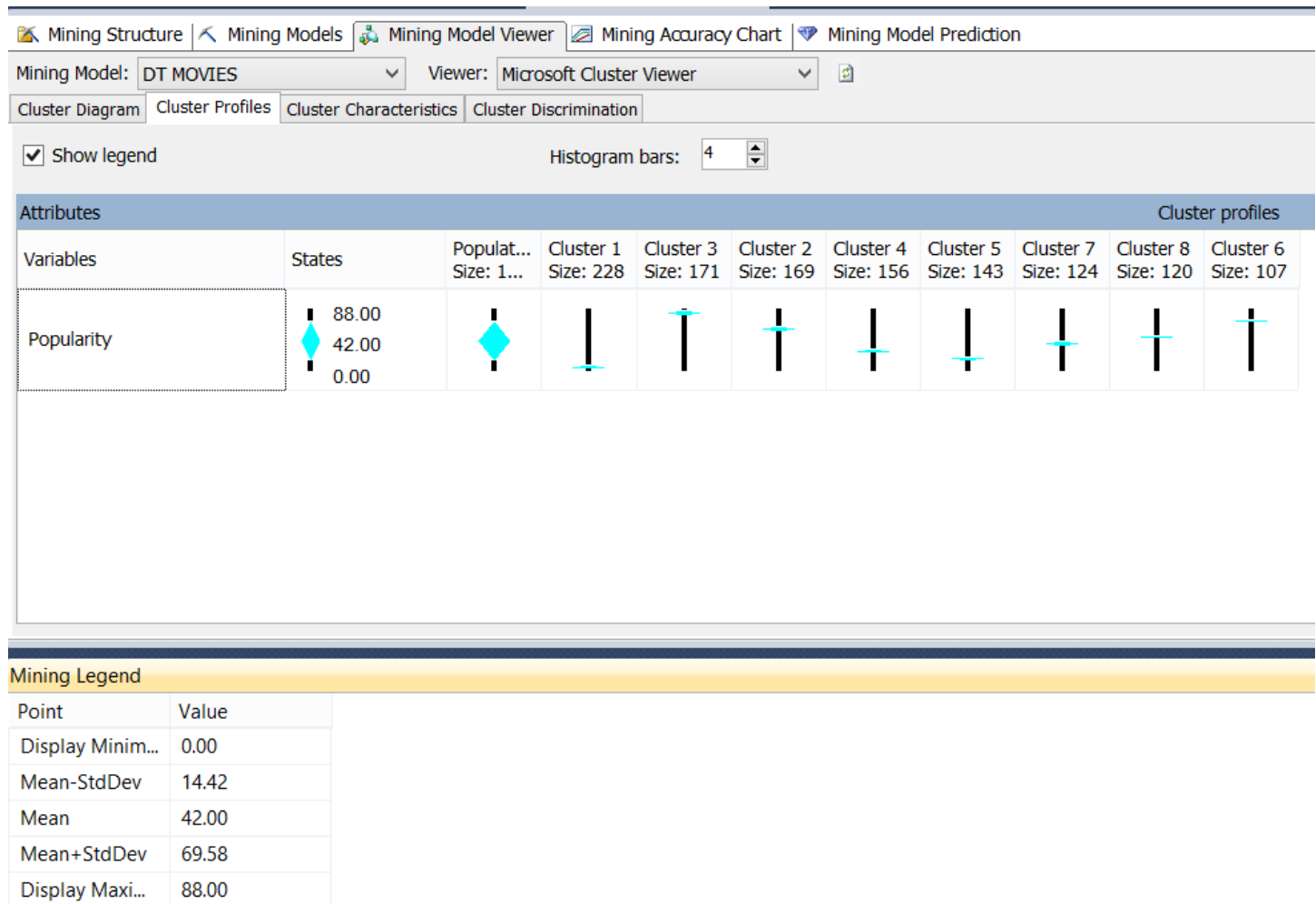## *Mining model for the clustering algorithm*



The Mining Model elaborates the relation between the nested table "FT_MASTER" and case table "DT_MOVIES". The attribute Popularity serves to be the Predict Column and Movie_id attribute being the Primary Key. The Data's are clustered using both the attributes and the calculation happens with Predict.

## *Screenshots of working Cluster Algorithms*

Given below is the result of the clustering algorithm we have used on our movie database for the above scenario. The links between the clusters shows how strongly they are correlated based on the distance. Darker the color, stronger is the link. The Cluster with Strong Color pattern indicates the more insight of data's grouped with in the desired range. The Cluster 1 is the darkest of all which gives an overview depicting more rates are given in this desired range.
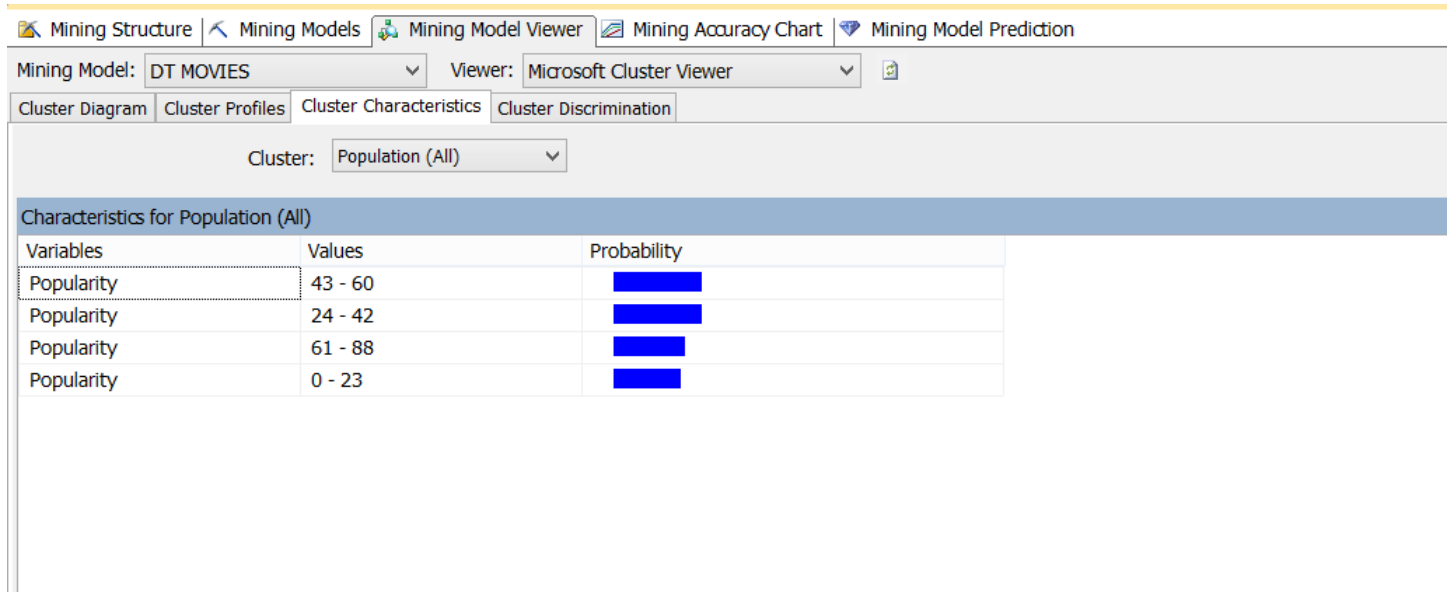
## *Cluster Profile of the Algorithm is used*



Here we see the averages and ranges of the individual clusters. The average popularity for the movies is 42 and if you will notice the most clusters have a varying quality to them for example Cluster 2 is very low, Cluster 3 is very low and Cluster stands to be entire variant to other two being very high. The Histograms formed classifies the number of data's in each buckets the each cluster.

## Cluster Characteristics



The Popularity ranges have been well split among the rated values, the values and chart are well determined in the decreasing manner. If you see for example more number of users have rated movies between 43 -60 and least was rated between 0 to 23. While the first two ranges all have approximately the same probability the low range has a much lower probability. This is good since you do not want to see a lot of movies with such a low rating

## 5.   Interesting Mining Results

The cluster are created with proper weightages, the buckets are split in a fashion that the popularities are spread across the buckets properly. For Example the clustered one is sized accordingly with average and has the maxium size 1218. The Most interesting of all was grouping of clusters, there were 6 clusters formed with right set of data distribution.

## 6.   Challenges

- Data Mining operations can be performed only the cleanse data, so getting the cleanse data was one among the challenges.

- Data mining operations have to be performed on fact and dimension tables, forming multiple relations and hierarchy was really challenging.

- Forming the clusters was the biggest challenge among all these, we have to make sure the entity relation is perfectly right and data are fully cleansed.

- Finding the best relation between the nested table and the case table was really hard.