# The Influence of Chance on Data Mining Results

*By: Ammar Y. Elnour*          *Supervisor : Prof. Joseph Davis*
School of Information Technology

FACULTY OF ENGINEERING & INFORMATION TECHNOLOGIES

> **A man with one watch knows what time it is; a man with two watches is never quite sure**   (An old saying)

## 1. Objectives

- To develop a generalized conceptual model to asses the influence of chance on the results of data mining. The concept will be based on the theory of statistical inference and uncertainty when applied on scale-free distributions.

- To develop testing procedures to test the method using nonparametric and bootstrapping methods and techniques such as hypotheses testing, statistical estimation and Inference.

- To validate the methods using experiments on large datasets.

## 2. Introduction

The rapid developments in information technologies have rendered it easy to gather and electronically capture a huge volume of data in different fields. These volume of data has the potential to produce useful information and new knowledge if appropriate mining techniques are used, however a good proportion of the results are redundant, obvious or statistically invalid. This research addresses the validity of the results which is still a rich area for research. Statistical techniques alone are not capable enough to address mining issues [1], but they can play significant role in assessing the results.

## 3. Chance-setup in Data Mining

### 3.1 What is chance setup

The chance is a dispositional property of one or more trials, experiments, or observations. This property relates to the long run frequency of an outcome on the trials [2].

### 3.2 The influence of chance on data mining

Data mining datasets are observational data. They are not controlled by any method of statistical experimental techniques which means that data mining researchers have no control over the source of the data, which leads to the lack of the possibility for random assignment to different group of data [3]. Datasets used in data mining are relatively large.
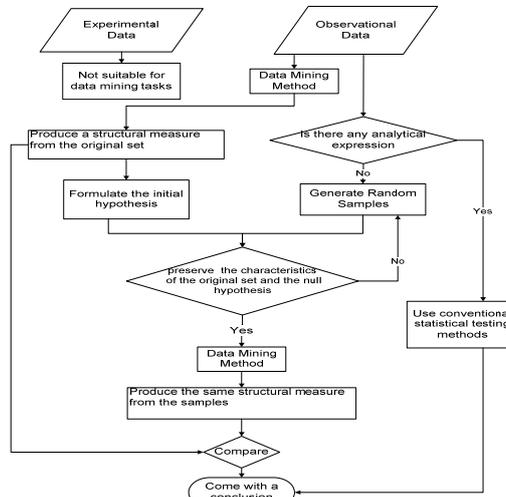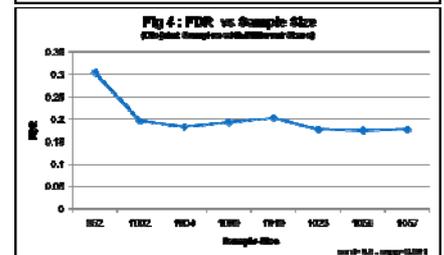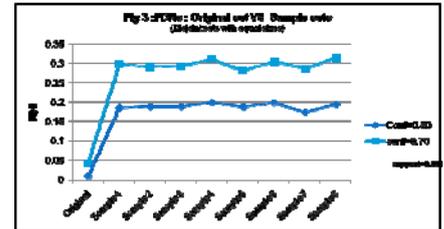


**Fig1** : Testing the validity of data mining results

Large data sets are more likely to be less clean which is an initial requirement for most statistical methods of analysis. The identically and independently distribution of data items is also one of the major assumptions that most statistical methods based on. However, large data sets are unlikely to fulfil this assumption. Data sets used in data mining tasks have their own characteristics where their chance set-up exists and affects later manipulations or processes. Datasets used in data mining tasks may be the result of mixture systems [4], where some hidden factors, unrecorded causes or latent variables may be the cause of part of the discovered associations and patterns. Considering many factors that may trigger chance set-up in observational datasets used on data mining tasks and the fact that data mining results offer candidate theories which need further validations using statistical approaches [6]. It is reasonable to conjecture that chance can influence the results of data mining methods.

## 4. Research Methods and Tools

Statistical inference methods are the main tools in our research. Bootstrap and permutation methods will be used to simulate and approximate the distribution of the original datasets. Nonparametric bootstrap hypotheses testing methods will also be used to test the validity of discovered knowledge and hence to quantify the influence of chance on the results. *Figure 1* shows the inferential model that we are proposing to test the validity of discovered knowledge from a data mining task.




## 5. Experimental Results

### 5.1 The influence of chance on association rule mining

Chance setup has significant effect on the generated results by association rule mining (ARM) .
We used the method of false discovery rate (FDR) on ARM [5] and showed by experiments that chance has noticeable impact on the results especially when low ARM is used. Figure 2,3 and 4 illustrate some of the findings.

### 5.2 Testing the validity of specific association rules

The FDR method can be used to validate a bulk set of rules under a specific range of supports and rule-confidences.
We introduce a method that use nonparametric bootstrap testing to test the validity of a single rule. Given a rule (R) of the form R: A→B under a fixed min support s and rule-confidence c. Our null hypothesis is that : The statistics of the rule-confidence denoted by *conf(R)* is equal to $c(\theta)$
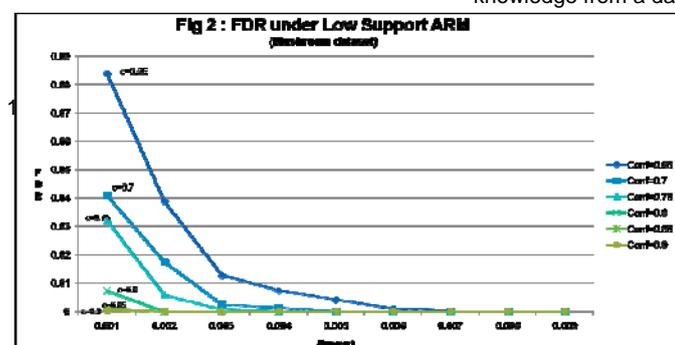i.e. $H_0 : conf(R_i) = c(\theta)$

We calculate the empirical *p-value* from the simulated samples $p = \frac{1 + \#\{t^* \geq t\}}{N + 1}$
Applying our method on the Mushroom dataset from the UCI we found that some rules which appear to be interesting are statistically invalid.

### 5.3 Further work:

We are working on a method to generate random samples that reflect the chance setup and suitable for assessing the results of a general data mining task.

## 6. References:

[1] M. David J. Hand and P. Smyth, Principles of Data Mining. MIT Press, 2001.
[2] Hacking, Logic of Statistical Inference. Cambridge University Press, 1965.
[3] D. Hand, "Mining the past to determine the future: Problems and possibilities," International Journal of Forecasting, October 2008.
[4] C. Glymour, D. Madigan, D. Pregibon, P. Smyth, and U. Fayyad, "Statistical themes and lessons for data mining," Data Mining and Knowledge Discovery, vol. 1, no. 1, pp.11–28, 1997.
[5] Zhang and B. Padmanabhan, "Using randomization to determine false discovery rate for rule discovery," in Proceedings of the Fourteenth Workshop On Information Technologies and Systems 11,12 December 2004. WITS, 2004, pp. 140–145.
[6] G. Pandey, S. Chawla, S. Poon, B. Arunasalam, and J. G. Davis, "Association rules network: Definition and applications," Stat. Anal. Data Min., vol. 1, no. 4, pp. 260–279, 2009.