



HAL
open science

Advanced Technology and Social Media Influence on Research, Industry and Community

Reda Alhajj

► **To cite this version:**

Reda Alhajj. Advanced Technology and Social Media Influence on Research, Industry and Community. 6th IFIP International Conference on Computational Intelligence and Its Applications (CIIA), May 2018, Oran, Algeria. pp.1-9, 10.1007/978-3-319-89743-1_1 . hal-01913884

HAL Id: hal-01913884

<https://inria.hal.science/hal-01913884>

Submitted on 6 Nov 2018

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

Advanced Technology and Social Media Influence on Research, Industry and Community

Reda Alhajj

Dept. of Computer Science, University of Calgary, Calgary, Alberta, Canada
e-mail: alhajj@ucalgary.ca

Abstract The rapid development in technology and social media has gradually shifted the focus in research, industry and community from traditional into dynamic environments where creativity and innovation dominate various aspects of the daily life. This facilitated the automated collection and storage of huge amount of data which is necessary for effective decision making. Indeed, the value of data is increasingly realized and there is a tremendous need for effective techniques to maintain and handle the collected data starting from storage to processing and analysis leading to knowledge discovery. This chapter will cite our accomplished works which focus on techniques and structures which could maximize the benefit from data beyond what is traditionally supported. In the listed published work, we emphasized data intensive domains which require developing and utilizing advance computational techniques for informative discoveries. We described some of our accomplishments, ongoing research and future research plans. The notion of big data has been addressed to show how it is possible to process incrementally available big data using limited computing resources. The benefit of various data mining and network modeling mechanisms for data analysis and prediction has been addressed with emphasize on some practical applications ranging from forums and reviews to social media as effective means for communication, sharing and discussion leading to collaborative decision making and shaping of future plans.

Keywords social media, social networks, data analysis, big data, frequent pattern mining, clustering, bioinformatics.

1 Introduction

Data is a major resource for decision making. Its value and importance has never been ignored since the existence of mankind on earth. It has been collected, stored and maintained using a wide variety of affordable means ranging from primitive to advanced. Indeed, collecting, storing and maintaining data was a cumbersome task in the past, mainly prior to the development of various technologies that gradually helped humans in handling data. However, the recent development in technology rapidly influenced data collection, storage and maintenance. For instance, sensors are becoming popular in all aspects of the daily life; they have been installed in almost every indoor and outdoor equipment. They are widely available and equipped with wireless communication skills which allow them to feed huge amounts of data

that should be captured, stored, cleaned, and processed for knowledge discovery as main ingredient of effective decision making.

In the past, humans used computing devices in a limited way. Database management systems were developed to facilitate flexibility in storing and retrieving data. Making sense of data was left to domain experts who are expected to retrieve and study data related to a specific problem in a way to draw some conclusions which may guide the decision-making process. Automating the knowledge discovery process was better realized towards the end of the 20th century when various machine learning and data mining techniques were developed and put in practice to serve a variety of application domains including business, health, security, etc.

To cope with the new era, researchers, developers and practitioners realized the need to develop new techniques and technologies capable handling growing volumes of data captured incrementally from heterogeneous sources. In other words, growth in volume and types of data expected to be processed suddenly witnessed a boom. Social networks and social media platforms are gaining increased popularity and are generating tremendous amounts of data. Surveillance devices are available almost everywhere. Even traditional archives are digitized. Consequently, storage media and techniques which were previously accepted as sufficient are no more capable of handling new needs. For instance, hard drives of personal computers were only couple megabytes in capacity when they were initially manufactured with less than one megabyte of main memory. People were happily competing to get the honor of owning and using such devices. It may be impossible to image using same computing platform in the current era where gigabytes of storage are no more sufficient. In fact, computing resources have improved rapidly to partially meet human needs but will never be satisfactory. Therefore, researchers and developers are always seeking new technologies and techniques, and hence conducting and advancing research will continue to attract more attention and investment.

Explicitly speaking, data volume, characteristics and associated expectations may be described as a moving target. This necessitates the availability of enough room for storage and sophisticated techniques for processing. People will continue to collect more data as time passes, but they will never afford to increase their computing power to handle their data effectively. Thus, the need for algorithms and techniques that can depend only on limited computing resources to deal with various aspects of data from dynamicity to volume, among others. Along this direction, we contributed various techniques and algorithms that could successfully satisfy a variety of applications which require handling large volumes of dynamic and stream data. These techniques are described in our published papers listed in the references at the end of this paper. Scalability is the main aspect considered by our techniques, including frequent pattern mining, clustering, network analysis, finding repeating patterns in long sequences, etc.

2 Partial Mentioning of Our Achievements

Our completed and ongoing research addresses various aspects of data from definition to construction to manipulation and analysis leading to knowledge discovery for decision making. Our initial contributions focused on traditional aspects related to handling and manipulation of data which were popular during the last two decades of the 20th century. We then gradually moved onto more advanced techniques which we realized as necessities since 1990s. These techniques, include, network analysis, data mining and machine learning techniques which have tremendously and visibly served various applications. We also realized scalability as a serious need especially in the current era of big data. We developed advanced techniques and adapt them to various domains, including:

- Bioinformatics and Health informatics
- Data partitioning and allocation
- Homeland security and terror/criminal network analysis
- Financial data analysis: from stock market to FOREX to fraud detection
- Web/network data analysis: from structure to content to usage
- Social media analysis and opinion mining including spam detection
- Recommendation and customer behavior analysis
- Network representation is a powerful mechanism for modeling many-to-many relationships.

A network consists of a set of nodes corresponding to the entities in the application domain and a set of links representing certain types of relationships between the entities. On the other hand, data mining includes a set of powerful techniques for studying the relationships/connections between various objects. Further, data mining may be used in the network construction phase. To construct a network data may be first analyzed using techniques like frequent pattern mining or clustering. Once a network is constructed, it can be analyzed for knowledge discovery.

Most of the traditional approaches for frequent pattern mining assume unlimited main memory which is not realistic. Therefore, scalability is a major concern when it comes to practical applications where data streams dynamically and available in large volume. To tackle these problems, we developed a novel approach which satisfies the following:

- The ability to mine in a bounded amount of memory space that may vary based on task priority. Thus, it is possible to mine using common PC.
- Improve external data access and make the mining process more I/O conscious.
- Introduce a specialized mining task aware memory manager for both RAM and the external memory

We build a tree structure namely, Frequent-Pattern (FP) tree, which summarizes the given data and allows for effective discovery of frequent patterns. Each branch represents at least one transaction. We build the tree from left to right and from top to bottom as shown in Figure 1.

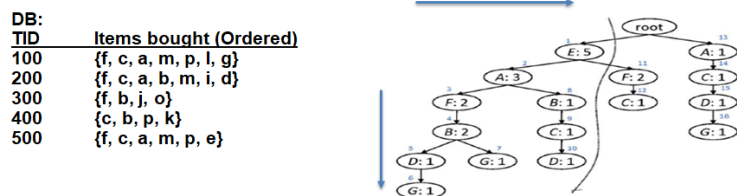


Fig. 1 Construction of FP-tree top-down and left-to-right

This way, we can store on the disk left side of the tree as it grows to the right. Therefore, our upper bound is the size of free disk space rather than the available memory.

NetDriller : A Powerful Social Network Analysis Tool*

- **Social Network Analysis (SNA)** is a technique first used in sociology.
- Recently computer scientists have realized that this model is general enough to be applied to **any domain** where the entities and their interconnections can be separated into **actors** and their **links**, respectively.
- **Data Mining** techniques can strengthen SNA

age	work class	education	Marital status	occupation	relationship	race	sex	Hours/week	native country
39	Private	Bachelors	Never-married	Adm-clerical	Not-in-family	White	Male	40	US
50	Self-emp-not-inc	Bachelors	Married-coupled	Exec-managerial	Partner	White	Male	45	Canada
30	Self-emp-not-inc	HS grad	Married-coupled	Exec-managerial	Partner	White	Male	45	US
52	State-gov	Bachelors	Married-coupled	Prof-specialty	Partner	Black	Male	40	US
25	Self-emp-not-inc	HS grad	Never-married	Farmer-fishing	Own-child	White	Male	35	US
41	Self-emp-not-inc	Masters	Divorced	Exec-managerial	Partner	White	Female	45	US

Raw Dataset: People and their attributes

1 Network Construction

Social Network: Based on community detection

2 Searching in the Network

Example 1: Find individuals who could monitor the information flow in an organization better than most others.

Example 2: Find individuals who have best picture of what is happening in the network as a whole.

- **Closeness** centrality reveals how long it takes information to spread from one individual to others in the network. High scoring individuals in Closeness have the shortest paths to all others in the network.
- **Betweenness** centrality indicates the extent that an individual is a broker of indirect connections among all other in a network. Someone with high Betweenness could be thought of as a gatekeeper of information flow. People that occur on many shortest paths among other People have highest Betweenness value.
- **Degree** centrality indicates the extent that an individual send or receive information to the neighbors.
- **Eigenvector** centrality calculates the principle eigenvector of the network. A node is central to the extent that its neighbors are central.

Fuzzy Sets: Based on multi-objective GA optimization

Fuzzy Query Result: Color hue shows DotM

* ICDM 2011 IEEE International Conference on Data Mining

<http://cpsc.ucalgary.ca/~nkoochak/NetDriller/>

Fig 2. Basic Characteristics of NetDriller.

To facilitate effective data investigation and analysis, we build our own tool, namely NetDriller which is capable of analyzing raw data to derive a network. Then various network analysis techniques could be applied on the network to identify actors which may reveal some important aspects related to the analyzed network, like most knowledgeable employee, most dangerous criminal, least performing student, best team to undertake next project, etc. The basics of NetDriller are summarized in Figure 2.

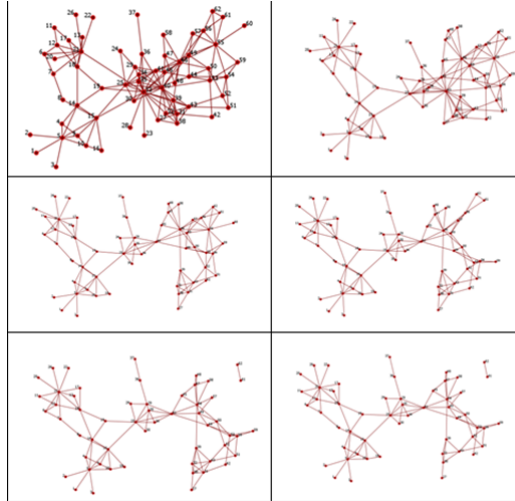


Fig 3. Sep. 11 network changes after excluding each level nodes for eigenvector centrality measure.

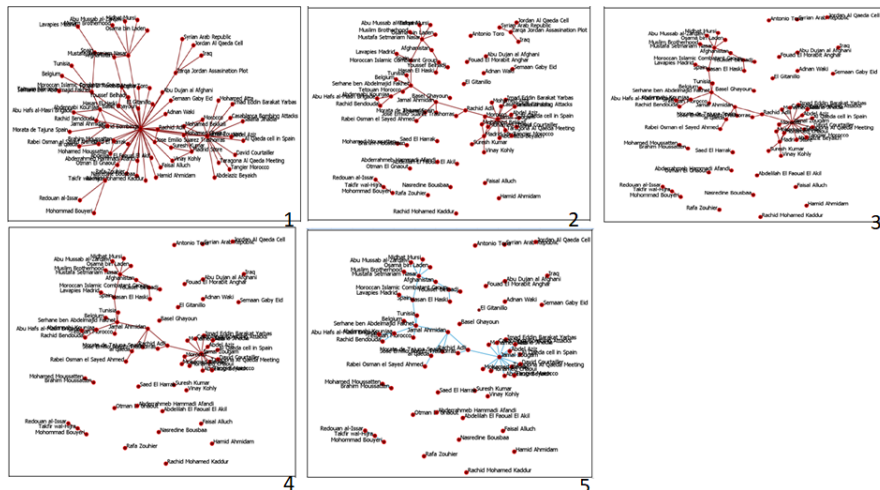


Fig 4. Madrid network changes after excluding each level nodes for eigenvector centrality measure

We utilized NetDriller to analyze September 11 terror network. It is surprising to realize that those who planned for the attacks considered all difficulties they could face. In other words, the network continues to be connected after removing terrorists who were identified as leaders down to level 6. The same is not true when the network of Madrid attacks was analyzed. The latter network became disperse only after removing second level leaders as shown in Figure 4.

Genes	Keywords
PCAP	regulate
HPC5	interact
MAD1L1	activate
HPC4	suppress
HIP1	prognostic
MSR1	biomarker
KLF6	network
PTEN	prognosis
MXI1	elucidate mechanism
CD82	prostate pathway
BRCA2	
CDH1	
ZFH3	
ELAC2	
HPCQTL19	
HPC3	
CHEK2	
HPC6	
AR	

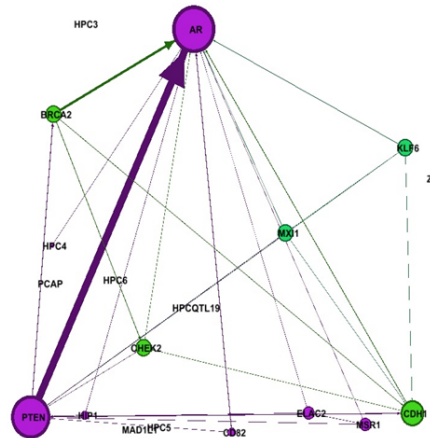


Fig 5. A list of genes and part of the gene-gene network related to prostate cancer.

NetDriller was also employed using gene expression data related to prostate cancer to identify proteins attributed to the disease. The main result reported by Net-Driller is shown in Figure 5.

Finally, next are some of our ongoing and planned research activities based on the promising results reported in our already published papers. First, in bioinformatics we are tracking disease evolution: spatial and temporal aspects, drug repositioning, etc. Second, in health Informatics we are working on patient monitoring, referral optimization and prediction, etc. Third, we demonstrated the applicability and effectiveness of sequence analysis and prediction for various domains, including financial (e.g., stock, forex), weather, traffic, energy, etc. Finally, other domains and applications considered by our research recommendation, sentiment analysis, opinion mining, spam detection, homeland security, close monitoring and analysis for early warning, etc.

To sum up, our research efforts described in our papers published in the literature and listed in the bibliography illustrate how data mining and network analysis are powerful techniques for data analysis. Further, it is possible to analyze huge amounts of data using limited computing resources, and to develop integrated solutions by combining various aspects leading to robust framework. We have succeeded in developing some techniques from scratch and we also expanded some existing techniques to produce working solutions for our industrial and academic partners. We could help in sophisticated data analysis to maximize knowledge discovery for informative decision making.

Bibliography

1. U. Manber and Myers, G., "Suffix Arrays: A New Method for On-Line String Searches." In Proceedings of the first annual ACM-SIAM symposium on Discrete algorithms (1990), pp. 319-327.
2. G. Cormode and Hadjieleftheriou, M., "Methods for finding frequent items in data streams," The VLDB Journal, (2009), unpaginated, doi: 10.1007/s00778-009-0172-z
3. R.S. Boyer, Moore, J., "A fast majority vote algorithm," Technical Report ICSCA-CMP-32, Institute for Computer Science, University of Texas (1981)
4. E. Demaine, López-Ortiz, A., Munro, J.I., "Frequency estimation of internet packet streams with limited space," In: European Symposium on Algorithms (ESA) (2002)
5. R. Karp, Papadimitriou, C., Shenker, S., "A simple algorithm for finding frequent elements in sets and bags," ACM Trans. Database Syst., (2003), 28, pp. 51–55
6. G. Manku, Motwani, R., "Approximate frequency counts over data streams," In: International Conference on Very Large Data Bases, (2002), pp. 346–357
7. A. Metwally, Agrawal, D., Abbadi, A.E., "Efficient computation of frequent and top-k elements in data streams," In: International Conference on Database Theory (2005)
8. M. Greenwald, Khanna, S., "Space-efficient online computation of quantile summaries," In: ACM SIGMOD International Conference on Management of Data (2001)
9. N. Bandi, Metwally, A., Agrawal, D., Abbadi, A.E., "Fast data stream algorithms using associative memories," In: ACM SIGMOD International Conference on Management of Data (2007)
10. N. Alon, Matias, Y., Szegedy, M., "The space complexity of approximating the frequency moments," In: ACM Symposium on Theory of Computing, pp. 20–29, 1996. Journal version in Journal of Computer and System Sciences, (1999), 58, pp. 137–147
11. G. Cormode, Muthukrishnan, S., "An improved data stream summary: the count-min sketch and its applications," J. Algorithms, (2005), 55(1), 58–75
12. K. F. Xylogiannopoulos, P. Karamelas, R. Alhajj: Real Time Early Warning DDoS Attack Detection. International Journal of Cyber Warfare and Terrorism, Volume 7, Issue 3, 2017.
13. S. Üçer, Y. Koçak, T. Ozyer and R. Alhajj, Social Network Analysis-based Classifier (SNAC): A Case Study on Time Course Gene Expression Data, Computer Methods and Programs in Biomedicine, 150(C):73-84, Oct. 2017.
14. A. Aksac, T. Ozyer and R. Alhajj, Complex Networks Driven Salient Region Detection based on Superpixel Segmentation, Pattern Recognition, Volume 66, June 2017, Pages 268–279.
15. G. Jurca, O. Addam, A. Aksac, S. Gao, T. Ozyer, D. Demetrick and R. Alhajj, "Integrating Text Mining, Data Mining, and Network Analysis for Identifying Genetic Breast Cancer Trends", BMC Research Notes, 9(1) · December 2016.
16. K. F. Xylogiannopoulos, P. Karamelas, R. Alhajj: Repeated Patterns Detection on Big Data Using Classification and Parallelism on LERP Reduced Suffix Arrays. Applied Intelligence, Vol.45, Issue 3, pp 567–597, 2016.
17. M. G. Ozsoy, F. Polat, and R. Alhajj, "Making Recommendations by Integrating Information from Multiple Social Networks", Applied Intelligence, July 2016, DOI: 10.1007/s10489-016-0803-1.
18. O. Addam, A. Chan, W. Hoang, R. Alhajj and J. Rokne. Foreign Exchange Data Crawling and Analysis for Knowledge Discovery Leading to Informative Decision Making. Knowledge-Based Systems. Volume 102, Pages 1–19, 2016.
19. A. Chen, A. Elhajj, S. Gao S. Afra, A. Sarhan, A. Kassem, and R. Alhajj, Approximating the Maximum Common Subgraph Isomorphism Problem with a Weighted Graph, Knowledge-Based Systems, Volume 85, pp.265–276, 2015.
20. O. shafiq, R. Alhajj and J. G. Rokne, On Personalizing Web Search using Social Network Analysis, Information Sciences, Vol. 314, 1 September 2015, Pages 55–76

21. K. F. Xylogiannopoulos, P. Karamelas, R. Alhajj: Analyzing very large time series using suffix arrays. *Appl. Intell.* 41(3): 941-955 (2014)
22. A. Rahmani, A. Chen, A. Sarhan, J. Jida, M. Rifaie and R. Alhajj, Social Media Analysis and Summarization for Opinion Mining: A Business Case Study, *Social Network Analysis and Mining*, (2014) 4:171.
23. K. F. Xylogiannopoulos, P. Karamelas, R. Alhajj, "Experimental Analysis on the Normality of pi, e, phi and square root of 2 Using Advanced Data Mining Techniques", *Experimental Mathematics*, Vol.23, No.2, pp.105-128, 2014
24. A. Rahmani, S. Afra, O. Zarour, O. Addam, R. Aljomai, N. Koochakzadeh, K. Kianmehr, R. Alhajj, Graph-based Approach for Outlier Detection in Sequential Data and Its Application on Stock Market and Weather Data, *Knowledge-Based Systems*, Vol.61, pp.89-97, May 2014.
25. W. Almansoori, O. Addam, O. Zarour, A. Sarhan, M. Elzohbi, M. Kaya, J. Rokne, R. Alhajj, The Power of Social Network Construction and Analysis for Knowledge Discovery in the Medical Referral Process, *Journal of Organizational Computing and Electronic Commerce*, Volume 24, Issue 2-3, 2014.
26. A. Qabaja, M. Alshalalfa, E. Alanazi and R. Alhajj, Prediction of Novel Drug Indications Using a Network Driven Biological Data Prioritization and Integration, *Journal of Cheminformatics*, 7;6(1):1, Jan 2014.
27. P. Peng, O. Addam, M. Elzohbi, S. Özyer, A. Elhajj, S. Gao, Y. Liu, T. Özyer, M. Kaya, M. Ridley, J. Rokne, R. Alhajj Analyzing Alternative Clustering Solutions by Employing Multi-Objective Genetic Algorithm and Conducting Experiments on Cancer Data, *Knowledge-Based Systems*, 56: 108-122, 2014.
28. M. Kaya and R. Alhajj, Development of Multidimensional Academic Information Networks with a Novel Data Cube Based Modeling Method, *Information Sciences*, 265: 211-224 (2014)
29. F. Rasheed and R. Alhajj, "A Framework for Periodic Outlier Pattern Detection in Time Series", *IEEE Transactions on Systems, Man, and Cybernetics*, 2014 May;44(5):569-82.
30. P. Polash Paul, M. Gavrilova and R. Alhajj, Decision Fusion for Multimodal Biometrics using Social Network Analysis, *IEEE Transactions on Systems, Man, and Cybernetics*, 44(11): 1522-1533, 2014.
31. J. Szeto, A. Lycett, X. Yi, S. Afra, A. Sarhan, K. F. Xylogiannopoulos, P. Karamelas and R. Alhajj, Integrating Data Mining Techniques into a User-Friendly Framework for Visualization of Health Indicators, *Health Informatics*, (2014) 3:63
32. M. Alshalalfa and R. Alhajj, "Integrating Protein Networks for Identifying Cooperative miRNA Activity in Disease Gene Signatures", *BMC Bioinformatics*, 2013;14 Suppl 12:S1. doi: 10.1186/1471-2015-14-S12-S1.
33. O. Öztürk, A. Aksaç, A. M. Elsheikh, T. Özyer and R. Alhajj, A Consistency-Based Feature Selection Method Allied with Linear SVMs for HIV-1 Protease Cleavage Site Prediction, *PLOS One*, 23;8(8):e63145. 2013
34. A. Guerbas, O. Addam, O. Zarour, M. Nagi, A. Elhajj, M. Ridley and R. Alhajj, "Effective Web Log Mining and Online Navigational Pattern Prediction", *Knowledge-Based Systems*, Volume 49, September 2013, Pages 50–62.
35. W. Almansoori, S. Gao, T.N. Jarada, A. M. Elsheikh, A. N. Murshed, J. Jida, R. Alhajj and J. Rokne, "Link prediction and classification in social networks and its application in healthcare and systems biology", *Network Modeling and Analysis in Health Informatics and Bioinformatics*, Volume 1, Issue 1-2, pp 27-36, June 2012.
36. M. Nagi, A. Elhajj, O. Addam, A. Qabaja, O. Zarour, T. Jarada, S. Gao, J. Jida, A. Murshed, I. Suleiman, T. Özyer, M. Ridley, R. Alhajj, Robust Framework for Recommending Restructuring of Websites by Analyzing Web Usage and Web Structure Data, *Journal of Business Intelligence and Data Mining*, 7(1/2): 4-20, 2012.
37. M. Adnan, M. Nagi, K. Kianmehr, M. Ridley, R. Alhajj, and J. Rokne, Promoting where, when and what?: An analysis of web logs by integrating data mining and social network

techniques to guide eCommerce business promotions, Social Networks Analysis and Mining, 2012.

38. F. Rasheed, M. Adnan and R. Alhadj, "Out-Of-Core Detection of Periodicity from Sequence Databases", Knowledge and Information Systems, Volume 36, Issue 1, pp 277-301, July 2013
39. M. Khabbaz, K. Kianmehr and R. Alhadj, "Employing Structural and Textual Feature Extraction for Semi-Structured Document Classification", IEEE Transactions on Systems, Man, and Cybernetics-C, 42 (6): 1566 - 1578, 2012
40. F. Rasheed and R. Alhadj, "Periodic Pattern Analysis of Non-uniformly Sampled Stock Market Data", Intelligent Data Analysis, (16(6): 993-1011 (2012)).
41. Gao, S., J. Zeng, A.M. ElSheikh, G. Naji, R. Alhadj, J. Rokne and D. Demetrick, "A Closer look at "social" boundary genes reveals knowledge to gene expression profiles," Current Protein & Peptide Science. 12(7):602-13, Nov 2011.
42. M. Adnan and R. Alhadj, "A Bounded and Adaptive Memory-Based Approach to Mine Frequent Patterns from Very Large Databases", IEEE Transactions on Systems, Man, and Cybernetics-B, Vol.41, No.1, pp.154-72, Feb. 2011.
43. F. Rasheed, M. Alshalalfa and R. Alhadj, "Efficient Periodicity Mining in Time Series Databases Using Suffix Trees," IEEE Transactions on Knowledge and Data Engineering, Vol.23 No.1, pp.79-94, Jan. 2011.
44. M. Alshalalfa, T. Özyer, R. Alhadj and J. Rokne, "Discovering Cancer Biomarkers: From DNA to Communities of Genes," Int. Journal of Networks and Virtual Organizations, Vol. 8, Nos. 1/2, pp.158-172, 2011.
45. F. Rasheed and R. Alhadj, "STNR: A Suffix Tree Based Noise Resilient Algorithm for Periodicity Detection in Time Series Databases," Applied Intelligence, Vol.32, No.3, pp.267-275, June 2010.
46. F. Rasheed, M. Alshalalfa and R. Alhadj, "Adaptive Machine Learning Technique for Periodicity Detection in Biological Sequence," Journal of Neural Systems, Vol.19, No.1, pp.11-24, February 2009.
47. M. Adnan and R. Alhadj, "DRFP-Tree: Disk-Resident Frequent Pattern Tree," Applied Intelligence, Vol.30, No.2, pp.84-97, April, 2009.
48. M. Kaya and R. Alhadj, "Multi-Objective Genetic Algorithms Based Automated Clustering for Fuzzy Association Rules Mining," Journal of Intelligent Information Systems, Vol.31, No.3, pp.243-264, December 2008.
49. M. Kaya and R. Alhadj, "Online Mining of Fuzzy Multidimensional Weighted Association Rules," Applied Intelligence, Vol.29, No.1, pp.13-34, August 2008.