

ІНФОРМАЦІЙНІ ТЕХНОЛОГІЇ

UDC 519.2

R. T. KASUMOVA

Azerbaijan National Academy of Sciences, Institute of Information Technology, Baku, Azerbaijan

ON INTELLECTUAL ANALYSIS OF DOMAIN NAME REGISTRATION DATA

Анотація: Робота присвячена формуванню бази знань системи доменних імен, що відповідають інтересам Азербайджанської Республіки. З цією метою в роботі розробляється сховище даних для обробки великого обсягу реєстраційних даних доменних імен, проводиться кластеризація цих даних і генеруються правила для видобутку нових знань.

Ключові слова: домен, DNS, адміністратор, реєстратор, реєстрант, кластеризація.

Abstract: The paper is devoted to formation of knowledge base of the domain names system, serving to interests of the Azerbaijan Republic. According to this purpose, the data warehouse is developed for processing great volume of registration data of domain names, clustering data is performed, and rules are generated for extraction of new knowledge.

Key words: domain, domain name system, administrator, registrar, registrant, clustering, categorical data.

Introduction

As a global telecommunications network of information and calculation resources, Internet creates a global information space, serves as a physical basis for World Wide Web and a variety of systems (protocols) of data transfer. Domain name system is used for addressing of requests in Internet. Domain is a spatial region of domain names and is characterized by independence of subdomain allocation, inclusion of information systems in domain structure, availability of special information systems (DNS-servers) containing data on domain names, allocated in the domain, and executes the function of organization of domain name space [1]. Domain name is an identifier of a domain and (or) information system, possesses a unique structure conditioned as: limited set of symbols, a name that identifies the domain, which contains domain name and as a voluntary part domain label or host name unique within the limits of the domain, which is contained in the domain name.

Domain name carries out the functions of identification, individualization and addressing [2]. Analyses demonstrate that, domain names are also used as means for conduction of unfit and unethical competition. One of the examples of unethical use of Internet, is use of famous trademarks, service trademarks, place of origins of commodities, as well as brand names in domain names [3].

Relevance

Currently, deficiencies existing in the field of domain name registration, absence of transparency in regards to domain name registration process, violation of domain name registration rules by the registrars, purchase and use of domain name (in irresponsible manner) with the purpose of its further sale, non-existence of a single policy against invaders of domain names (cyber squatters, phishing etc), as well as software, which allow to conduct an accurate analysis of registration data about domain names collected in DNS (Domain Name System).

Listed problems make conduction of scientific analysis of domain name registration data collected in DNS servers necessary. Considering the dynamics of increasing of number of domain names, it is possible to conduct processing and analysis of this information collected from thousands of domain names, obtain new knowledge, detect regularities and make necessary decisions.

As intellectual analysis of registration data can be the solutions reason of such issues as forecasting, making operative, effective and analytical decisions in domain field, definition of facts, evaluation of the real condition of domain market, research domain name monitoring problems etc.

Objective

The objective of the article is development of data storage for processing of a large volume of domain name registration data, conduct clustering of this data and generates rules for gaining new knowledge.

Tasks

Following tasks are formulated in research purpose:

- 1) Processing of domain name registration data in storage, clustering of this data using CLOPE algorithm;
- 2) Generation of rules using Magnum Opus v.5.4.1. Program for each cluster and decision making.

Problem Solution

Domain registration information include: domain name, registrar, name, address, admin-o, admin-c, organization, created, updated, free-date, phone, e-mail, nserver, type, source, paid-till etc [4] (pic. 1).

Private Person	azerbajan.ru	ru	RUCENTER-REG-RIPN	2000	2011	Russian	+7.095930	-	d@bugtraq.ru	-	active	San Francisco, United States
Minakuman Periasamy	azerbajan.com	com	ENOM, INC.	1997	2017	Malaysia	347482144	1.65688730	nicolov@eml.cc	-	active	London, United Kingdom
M.A. Stenzel	azerbajan.net	net	GODADDY.COM, INC.	1999	2017	US	-	-	-	-	inactive	San Antonio, United States
Role	azerbajan.org	org	Moniker Privacy Services	2000	2011	US	+1.954984	+1.954969	support@moniker.com	sale	passive	United States
Reserved under ICANN Re	azerbajan.info	info	Internet Corporation for Assi	2001	2006	US	+1.310823	+1.310823	res-dom@iana.org	-	inactive	-
Jong Won Hwang	azerbajan.biz	biz	-	2002	2011	KOREA, REPUBLIC	+82.55687	+82.55687	jongwonhwang@yahoo	sale	passiv	Los Angeles, United States
-	azerbajan.tv	tv	ENOM, INC.	2003	2011	-	-	-	-	-	active	Germany
-	azerbajan.su	su	RUCENTER-REG-FID	1999	2011	Russian	+7.495737	+7.495737	ru-ncc@nic.ru	-	passive	United States
Midrag Radosavljevic	azerbajan.cc	cc	MESH DIGITAL LIMITED	2007	2011	Beograd, Serbia	+381.6337	-	domains@inbox.com	-	passiv	San Diego, United States
Government of India	azerbajan.in	in	RESERVED NAME	2004	2009	India	+91.11243	+91.11243	amar@ernet.in	-	inactive	-
Shajeem Othayoth	azerbajan.me	me	-	2010	2011	India	+91.98953	-	othayoth.domains@gmai	sale	passiv	Frankfurt Am Main, Germany
Private Domain Registratio	azerbajan.ws	ws	.WS Registry	2003	2014	US	+1.760602	-	azerbajan.ws@privated	-	passive	Carlsbad, United States
Telnic-Master-Contact	azerbajan.tel	tel	Telnic Ltd	2009	2011	UNITED KINGDOM	+44.20746	+44.20746	support@registry.nic.tel	-	inactive	-
-	baku.com	com	NETWORK SOLUTIONS, LLC.	1996	2012	US	-	-	-	-	inactive	Herndon, United States
-	baku.net	net	ENOM, INC.	2001	2011	Karachi, PK	+92213002	+92213002	info@dotcorner.com	sale	passive	Frankfurt Am Main, Germany
M.A. Stenzel	baku.org	org	GODADDY.COM, INC.	2000	2011	US	+1.808878	-	contact@stenzel.org	-	inactive	-
Private Person	baku.ru	ru	RUCENTER-REG-RIPN	2000	2011	Russian	+7.095930	-	d@bugtraq.ru	-	active	San Francisco, United States
Moniker Privacy Services	baku.info	info	Moniker Privacy Services	2001	2012	US	+1.954984	+1.954969	BAKU.INFO@domainserv	-	active	Houston, United States
Moniker Privacy Services	baku.biz	biz	Moniker Privacy Services	2005	2011	US	+1.954984	+1.954969	BAKU.BIZ@domainservic	-	active	Houston, United States
-	baku.tv	tv	ENOM, INC.	2004	2011	-	-	-	-	sale	passive	Houston, United States
Private Person	baku.su	su	RUCENTER-REG-FID	2005	2010	Russian	+7.095123	+7.095123	taragir@genocide.ru	-	active	Houston, United States
Domain Manager	baku.in	in	Travel India	2005	2011	India	+91.91674	-	info@travelindia.org.in	-	inactive	-
-	baku.ws	ws	Wild West Domains, Inc.	2004	2011	US	+1.480624	-	dns@jomax.net	-	active	Erfurt, Germany
Marina V Zyryanova	baki.ru	ru	RUCENTER-REG-RIPN	2005	2011	Russian	+7.903793	-	uspeshanya@inbox.ru	-	active	Moscow, Russian Federation
-	baki.com	com	MONIKER ONLINE SERVICES,	2001	2011	Taiwan	+886.9175	-	domadm@mustneed.com	sale	passive	Toronto, Canada
Domain Manager	baki.tv	tv	DOTSTER, INC.	2007	2011	Azerbaijan	+99450217	-	domain@azintergroup.cc	-	active	Vancouver, United States
-	baki.net	net	NETWORK SOLUTIONS, LLC.	1998	2011	US	+1.310929	-	question@gmail.com	sale	passive	Frankfurt Am Main, Germany
Registration Private	baki.org	org	Domains by Proxy, Inc.	2005	2011	US	+1.480624	+1.480624	BAKU.ORG@domainsbypr	sale	passive	Frankfurt Am Main, Germany
Eugene V Petrov	baki.su	su	RUCENTER-REG-FID	2010	2011	Russian	+7.963344	-	etot_domen_prodayotsy	-	inactive	Moscow, Russian

Picture 1 – Domain registration data example

As domain registration data mainly consist of categorical data (which cannot be regulated in space), application of traditional algorithms for objects' clustering is ineffective. Clustering – is a fundamental data analysis and Data Mining task that groups together similar objects. On modern level, clustering is frequently used as the first step at data analysis [5, 6].

High dimensionality (thousand fields) and large volume (hundred thousand and millions of records) of data base tables, complexity of metrics definition for calculation of distance among categorical data, very low productivity at pair-wise comparison of distance between points (k-means) at each iteration procedure on large record arrays, and sometimes even inapplicability require application of scaled algorithms of categorical data clustering. In most algorithms, metrics based on Euclidian distance concept is used as objects' proximity

$$\text{metrics } d(x, y) = \sqrt{\sum_{k=1}^n (x_k - y_k)^2}, \text{ where } x = (x_1, x_2, \dots, x_n), y = (y_1, y_2, \dots, y_n), \text{ thresholds } - \lambda$$

are given for cluster setting, if $d(x, y) < \lambda$, then $x, y \in A_i, i = \overline{1, c}$. But given metrics is not always effective, as it obliges clusters to have a spherical form that is not inherent to them. Consequently, known k-means clustering algorithms cannot achieve a satisfying result.

Currently, a variety of clustering algorithms were proposed for working with categorical data. But, they do not always meet abovementioned requirements. LargeItem is considered as one of the effective algorithms, which is based on optimization of some global criteria.

CLOPE algorithm, proposed in 2002 by a group of Chinese scientists, allows clustering task solution of not only categorical, but also any transactional data [7]. It provides a higher productivity and better quality of clustering in comparison with LargeItem algorithms and many other hierarchical algorithms. The key is that all features of objects are measured in nominal scale. However, before launching CLOPE, data must be brought to normalized form. It can have a form of a binary matrix, as in associative rules, as well as being a biunivocal mapping between a set of unique objects of the table and a set of whole numbers. CLOPE is easily counted and interpreted. During its operation, algorithm saves a small amount of information on each cluster in RAM and requires a minimal number of data set scanning. This allows its application for clustering of huge volumes of categorical data. CLOPE automatically selects a quantity of clusters; moreover it is regulated by a single parameter – repulsion coefficient [8].

Thus in reviewed case, CLOPE is one of the effective algorithms, based on which is the idea of maximization of global criteria – cost function Profit (C), which increases the proximity of transactions in clusters through increasing the cluster histogram parameter (pic.2).

Tr.	Cluster	Prop. 1	Prop. 2	Prop. 3	Prop. 4	Prop. 5	Prop. 6	Prop. 7	Prop. 8	Prop. 9	Prop. 10	Prop. 11	Prop. 12	Prop. 13
Tr. 1	1	Private Pers	azerbaijan.ru	ru	RUCENTER-	2000	2011	Russian	+7.095930	-	d@bugtraq	-	active	San Fran
Tr. 2	13	Mikukumi	azerbaijan.com	com	ENOM, INC.	1997	2017	Malaysia	347482144	1.6568730	nicolov@em	-	active	London, Uni
Tr. 3	3	M.A. Stentz	azerbaijan.net	net	GODADDY.C	1999	2017	US	-	-	-	-	inactive	San Antonic
Tr. 4	3	Role	azerbaijan.org	org	Moniker Priv	2000	2011	US	+1.954984	+1.954969	support@im	sale	passive	United Stati
Tr. 5	7	Reserved u	azerbaijan.info	info	Internet Co	2001	2006	US	+1.310829	+1.310823	res-dom@ia	-	inactive	-
Tr. 6	14	Jong Won H	azerbaijan.biz	biz	-	2002	2011	KOREA, REP	+82.55687	+82.55687	jongwonhw	sale	passiv	Los Angeles
Tr. 7	3	-	azerbaijan.tv	tv	ENOM, INC.	2003	2011	-	-	-	-	-	active	Germany
Tr. 8	12	-	azerbaijan.su	su	RUCENTER-	1999	2011	Russian	+7.495737	+7.495737	ru-icc@nic	-	passive	United Stati
Tr. 9	11	Miodrag Rai	azerbaijan.cc	cc	MESH DIGIT	2007	2011	Beograd, Se	+381.6337	-	domains@in	-	passiv	San Diego, I
Tr. 10	15	Government	azerbaijan.in	in	RESERVED I	2004	2009	India	+91.11249	+91.11249	amar@eme	-	inactive	-
Tr. 11	14	Shajaeem Ot	azerbaijan.me	me	-	2010	2011	India	+91.98953	-	pthayoth@di	sale	passiv	Frankfurt Ai
Tr. 12	6	Private Dor	azerbaijan.ws	ws	JWS Registr	2003	2014	US	+1.760602	-	azerbaijan	-	passive	Carlsbad, U
Tr. 13	15	Tehnic@Hast	azerbaijan.tel	tel	Tehnic Ltd	2009	2011	UNITED KING	+44.20746	+44.20746	support@re	-	inactive	-
Tr. 14	3	-	baku.com	com	NETWORK S	1996	2012	US	-	-	-	-	inactive	Herdon, U
Tr. 15	9	-	baku.net	net	ENOM, INC.	2001	2011	Karachi, PK	+92213002	+92213002	info@dotco	sale	passive	Frankfurt Ai
Tr. 16	3	M.A. Stentz	baku.org	org	GODADDY.C	2000	2011	US	+1.808878	-	contact@sb	-	inactive	-
Tr. 17	1	Private Pers	baku.ru	ru	RUCENTER-	2000	2011	Russian	+7.095930	-	d@bugtraq	-	active	San Fran
Tr. 18	3	Moniker Priv	baku.info	info	Moniker Priv	2001	2012	US	+1.954984	+1.954969	BAKU.INFO	-	active	Houston, U
Tr. 19	3	Moniker Priv	baku.biz	biz	Moniker Priv	2005	2011	US	+1.954984	+1.954969	BAKU.BIZ@	-	active	Houston, U
Tr. 20	3	-	baku.tv	tv	ENOM, INC.	2004	2011	-	-	-	-	sale	passive	Houston, U
Tr. 21	1	Private Pers	baku.su	su	RUCENTER-	2005	2010	Russian	+7.095123	+7.095123	taragr@ger	-	active	Houston, U
Tr. 22	5	Domain Mar	baku.in	in	Travel India	2005	2011	India	+91.91674	-	info@travel	-	inactive	-
Tr. 23	3	-	baku.ws	ws	Wild West C	2004	2011	US	+1.480624	-	dru@jmax	-	active	Erlfurt, Gem
Tr. 24	1	Marino V Zy	baku.ru	ru	RUCENTER-	2005	2011	Russian	+7.903793	-	uspehamyc	-	active	Moscow, RU
Tr. 25	9	-	baku.com	com	MONIKER, O	2001	2011	Taiwan	+886.9175	-	domain@r	sale	passive	Toronto, Ca
Tr. 26	11	Domain Mar	baku.tv	tv	DOTSTER, I	2007	2011	Azerbaijan	+99450217	-	domain@bz	-	active	Vancouver,
Tr. 27	3	-	baku.net	net	NETWORK S	1998	2011	US	+1.310929	-	question@q	sale	passive	Frankfurt Ai
Tr. 28	3	Registrar	baku.org	org	Domains bv	2005	2011	US	+1.480624	+1.480624	BAKU.ORG@	sale	passive	Frankfurt Ai
Tr. 29	1	Eugene V Pr	baku.su	su	RUCENTER-	2000	2011	Russian	+7.963344	-	etot_dome	-	inactive	Moscow, RU

Picture 2 –Domain Registration Data Clustering

Formula for calculation of global criteria – cost function looks like following: for given transaction base $D = \{t_1, t_2, \dots, t_n\}$ and r is the repulsion coefficient, to find such a splitting of $C = \{C_1, C_2, \dots, C_k\}$ that

$$Profit(C, r) \rightarrow \max \quad (1)$$

Where r regulates the level of similarity of transactions within the cluster –bigger is r , larger is the final number of clusters.

Cost function formula:

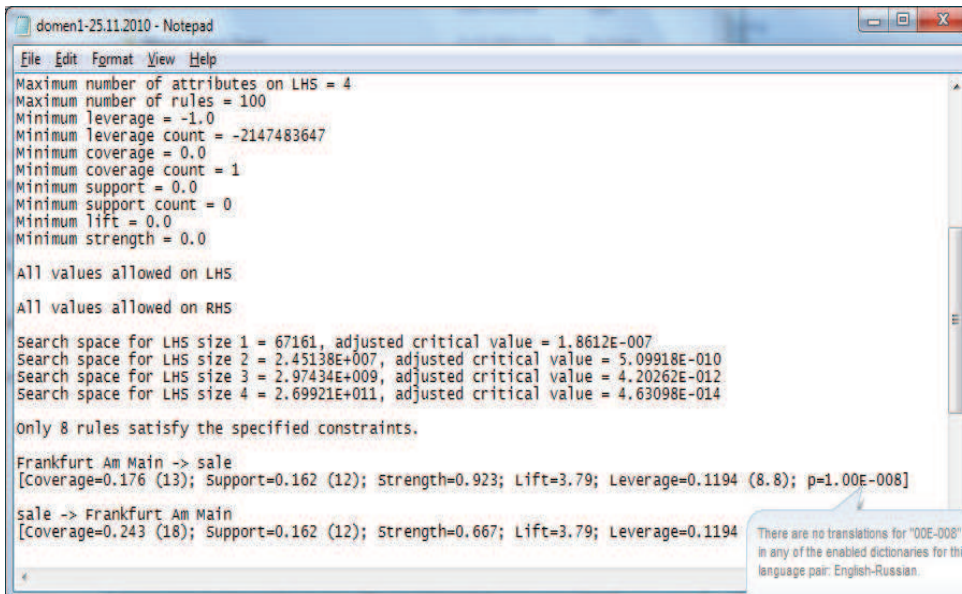
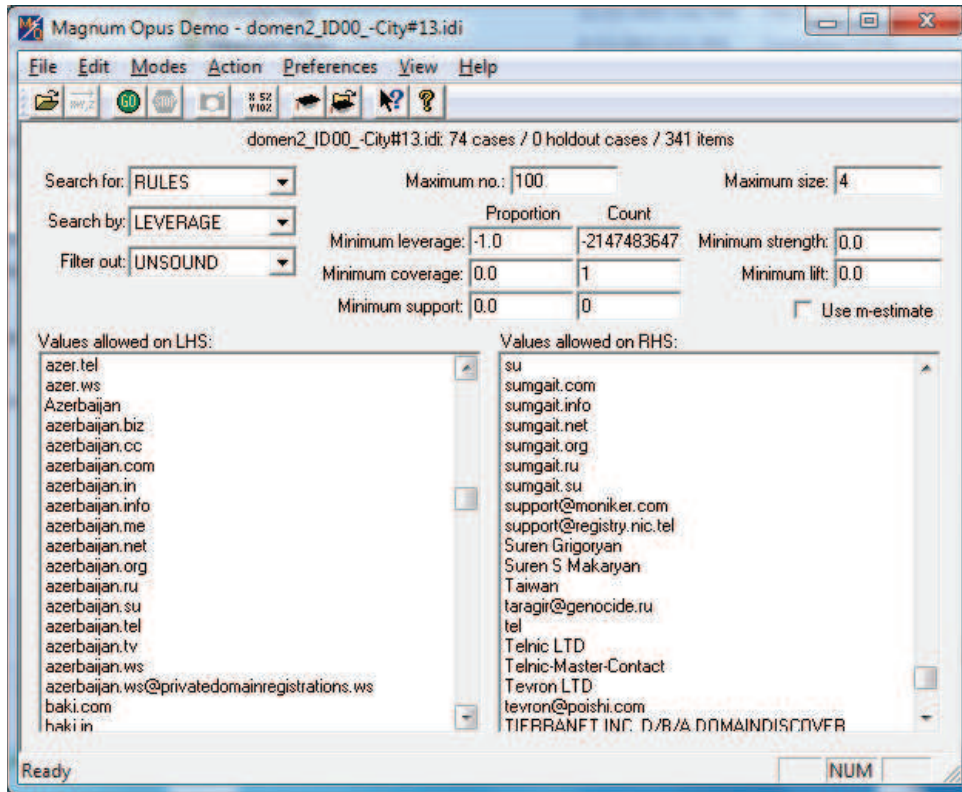
$$Profit(C, r) = \frac{\sum_{i=1}^k \frac{S(C_i)}{W(C_i)^r} \times |C_i|}{\sum_{i=1}^k |C_i|}, \quad (2)$$

where $S(C_i)$ – is the total number of entry of i objects in clusters C_i , $W(C_i)$ - is the quantity of unique objects C_i , $|C_i|$ - is the quantity of objects in i^{th} cluster, k – is the number of clusters, r - is the positive real number larger than 1.

After completion of clustering stage, we generate rules for each cluster using Magnum Opus v.5.4.1 program (<http://www.giwebb.com>), which is given case detects associations within the data framework of one cluster (pic. 3). Magnum Opus is easy-to-use and finds new associations in data. Associative rules are the results of the operation of his program.

In accordance in [9], associative rule is determined following way: let's assume that $I = \{i_1, i_2, \dots, i_n\}$ - a set of elements (binary attributes called elements). Any subset I is called an itemset of elements or simply itemset. Number of elements in the set is called its length. Let's assume that D - is a set of transactions, where each transaction $T, T \subseteq I$ - represents a binary vector, where $T.i_k = 1$, if element $i_k, k = \overline{1, n}$ is present in the transaction, otherwise $T.i_k = 0$ ($T.i_k$ – is the value of k^{th} attribute). Let's assume that X is a set of elements from I . Let's say that, transaction $T, T \subseteq I$ contains X , i.e. $X \subseteq T$, if $\forall i_k \in X$, then $T.i_k = 1$. Associate rule is $X \rightarrow Y, [c, s]$ type implication, where $X \subset I, Y \subset I, X \cap Y = \emptyset, c$ (confidence) - is

the confidence of the rule, s (*support, significance*) - is the support of the rule. Support of the rule - s is the percentage of transactions from D , containing $X \cup Y$, i.e. $s(X \rightarrow Y) = s(X \cup Y)$; and confidence of the rule demonstrates that c percentage of transactions from D containing X also contains Y , i.e. $c(X \rightarrow Y) = s(X \cup Y) / s(X)$. Abovementioned definition is an actual definition of Boolean (binary) associative rule, as it is only reviewed if $i_k, k = \overline{1, n}$ element is present in the transaction or not. In fact, transactions of data bases usually contain numerical and categorical data types.



Picture 3 – Rules generated for each cluster using Magnum Opus v.5.4.1 program

For example, experiments conducted with the domain name system knowledge base, serving the interests of the Republic of Azerbaijan, demonstrated that the majority of generic high level domains – gTLDs with the names of geographical locations of the country were purchased by foreign citizens living outside our republic in 1997-2001 years.

Conclusion

Today, undoubtedly, success of any research field directly depends on its capability of extracting and creating knowledge. With expansion of the commercial use of Internet capabilities, given sphere of human activity is the maximally demanded activity field for gaining profit as for right holders as well as for perpetrators. Usage of such inseparable attribute of presence in the Web as domain name opened new capabilities for lawful, as well as illegal use of such habitual commercial instruments, such as trade mark and brand name. For this reason, detection and replenishment of knowledge in developed domain name system knowledge base can become one of the main research directions in struggle against problems in registration qualification and use of domain names.

Reference

1. Серго А.Г. Доменные имена. М.: «Бестселлер», 2006, с. 368.
2. Милютин З. Ю. Правовой статус доменного имени //Патенты и лицензии. – 2005. № 6. - с. 19-23.
3. Калятин В. О. Проблема конфликта прав на доменные имена с правами на иные средства индивидуализации. // Юридический мир. – 2001. № 5. - с. 19-28
4. Венедрухин А.А. Доменные войны. СПб.: Питер, 2009, с. 224.
5. Дюк В., Самойленко А. Data Mining: Учебный курс. – СПб: Питер, 2001, с. 368.
6. Алгулиев Р.М., Касумова Р.Т., Алекперова И.Я. О современных концепциях, поддерживающих принятие решений //Известия Национальной Академии Наук Азербайджана. Серия физико-математических и технических наук. – 2005, №2, с. 70-75.
7. Yang, Y., Guan, H., You, J. CLOPE: A fast and Effective Clustering Algorithm for Transactional Data /KDD '02 Proceedings of the eighth ACM SIGKDD international conference on Knowledge discovery and data mining. July 23-26, 2002, Edmonton, Alberta, Canada, pp. 682-687.
8. Wang, K., Xu, C., Liu, B. Clustering transactions using large items. In Proc. CIKM'99, Kansas, Missouri, 1999.
9. Agrawal R., Imielinski T., Swami A. Mining association rules between sets of items in large databases. / In Proceedings of the ACM SIGMOD Conference on Management of Data, Washington D.C., May 1993, pages 207-216.

Відомості про авторів

Касумова Рена Тофік кизи – голова відділення Азейбаржанської національної академії наук, інституту інформаційних технологій, вул. Агаєва, 9, Баку, Азейбаржан. (99412) 439-85-48, kasumova-rena@rambler.ru.