



Gartner¹

2011 .,

2.

¹

²

... // CIO, 3, 2011, (http://cio.bg/3715_intelligentnite_biznes_analizi_tendencii_i_perspektivi, 2.07.2011).

I.

1.1.

(Data Mining – DM)

(Massachusetts Institute of Technology – MIT).

()

”4.

2,4 . 5.

³ , , 2009, . 68.

⁴ Hand, D., . Mannila and . Smyth. Principles of Data Mining. MIT Press, 2001, <http://mitpress.mit.edu/books/chapters/026208290Xchap1.pdf>, (23.08.2011).

⁵ Internet World Stats, <http://www.internetworldstats.com/stats.htm>, (8.12.2012)

1,2

5,32

2012

6.

spider crawler.

(,).

DM

7.

Web mining (WM).

1996 .T WM

WM

Web mining „

DM

„8.

WM – „

„9.

⁶ , http://cio.bg/4380_zakonat_na_mur_e_v_sila_i_za_internet_resursite, (10.01.2012).

⁷ Etzioni, . The World Wide Web: quagmire or gold mine?// Communications of the ACM, no. 11, 1996, p . 65–68.

⁸ Markov, Z. and D. Larosed. Data Mining the Web Uncovering Patterns in Web Content, Structure, and Usage. New Jersey: John Wiley & Sons, 2007.

⁹ Cooley, R, B. Mobasher and J. Srivastava. Web Mining: Information and Pattern Discovery on the World Wide Web. (<http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.27.3042&rep=rep1&type=pdf>, 3.07.2011).

WM, ¹⁰

• Web content mining (WCM).
 ()

WCM

• Web structured mining (WSM).
 (HT L, XML)

• Web usage mining (WUM).
 WUM

¹⁰ Cooley, R. Mobasher, B. and Srivastave, J., Web Mining: Information and Pattern Discovery on the World Wide Web (<http://maya.cs.depaul.edu/classes/ect584/papers/cms-tai.pdf>, 9.12.2011). Markov, Z. and D. Larosed. Data Mining the Web Uncovering Patterns in Web Content, Structure, and Usage. New Jersey: John Wiley & Sons, 2007.

¹¹ Cross Industry Standard Process for Data Mining - CRISP-DM (http://en.wikipedia.org/wiki/Cross_Industry_Standard_Process_for_Data_Mining, 31.08.2011); Fayyad, U., G. Piatetsky-Shapiro and P. Smyth. From Data Mining to Knowledge Discovery in Databases (<http://www.kdnuggets.com/gpspubs/aimag-kdd-overview-1996-Fayyad.pdf>, 27.08.2011).

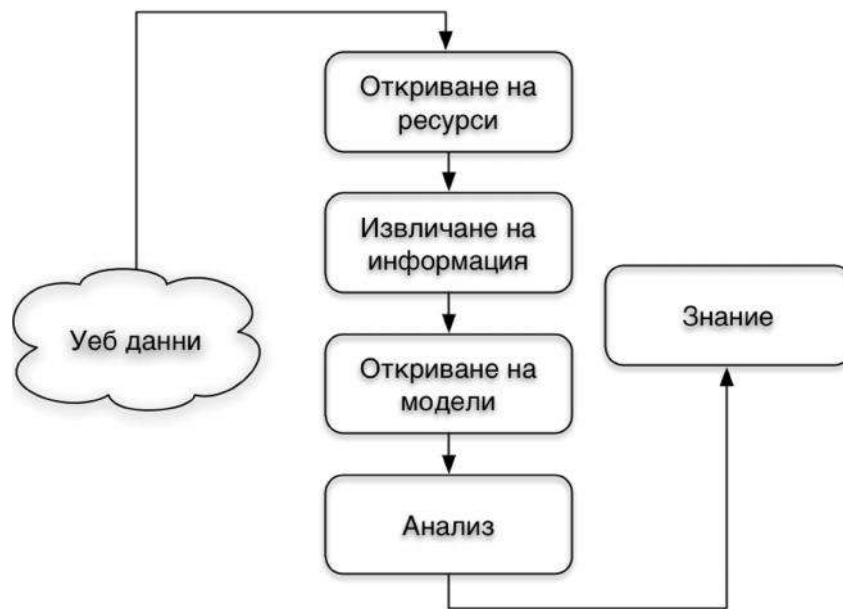
- 1.
- 2.
- 3.
- 4.
- 5.

DM

DM,

Web Mining

¹² (. 1):



. 1. WM

¹² Kosala, R. and H. Blockeel. Web Mining Research: A Survey (<http://facweb.cs.depaul.edu/mobasher/classes/ect584/papers/kosala.pdf>, 4.07.2011).

1.	-	.
	,	-
	,	-
2.	-	-
	.	-
	.	-
3.	,	.
	,	DM
	,	-
4.	-	-
	.	-
	,	-
	:	-
•	.	-
	.	-
	” - ”	-
	,	-
•	.	-
	-	-
	.	-
	,	e
	,	-
	,	-
•	,	.
	.	-
	.	-
•	.	.
	,	()
	.	-
	,	-
	.	-
•	.	.
	,	-

1-R Naive Bayes, ID3, C4.5 SVM.

1-R (1-rule)

1-R Naive Bayes.

1-R ID3, C4.5 „Gini”,

CART (Classification and Regression Tree).

ID3

13

0

0

C4.5

ID3.

CART

¹³ Shannon, Cl. Prediction and entropy of printed English. // The Bell System Technical Journal, 30 (http://www.princeton.edu/~wbialek/rome/refs/shannon_51.pdf, 12.01.2012).

Support Vector Machines (SVM),

Apriori.

BFS)¹⁴ . priori (breadth-first search

(Between-groups linkage)

(Nearest neighbor).

(K-Means Cluster).

k

k-

k

¹⁴

Web Mining.

1.2.

1. (), ()

Google. T

2. (stop words) –

3.

¹⁵ Fayyad, U., G. Piatetsky-Shapiro and P. Smyth. From Data Mining to Knowledge Discovery in Databases (<http://www.kdnuggets.com/gpspubs/aimag-kdd-overview-1996-Fayyad.pdf>, 5.07.2011).

1
4.

I
16

Ahonen	posode rules –		– –
Billsus Pazzani	– TFIDF – – Naive Bayes –		
Cohen	– Propositional rule based system – –		
Dumais	– TFIDF – – Naive Bayes – – Support Vector Machines	– –	
Feldman Dagan			
Feldman			
Frank	Naive Bayes –		

¹⁶ Kosala, R. and H. Blockeel. Web Mining Research: A Survey (<http://facweb.cs.depaul.edu/mobasher/classes/ect584/papers/kosala.pdf>, 4.07.2011).

Freitag McCallum			
Hofmann			
Honkela	Self-Organizing Maps –		
Junker			– –
Kargupta	– – –		– –
Nahm Mooney			– –
Nigam			
Scott Matwin	Rule based system -	– – – –	17
Soderland		– –	
Weiss			
Wiener	– –		
Witten		–	
Yang	– – –	– –	

HTML

. 2.

Craven	– Modified Naive Bayes – –		– – –
Crimmins	Unsupervised and supervised classification algorithms -	, URL,	– –
Fürnkranz			–
Joachims	– TFIDF – –		–
Muslea		,	
Shavlik Eliassi- Rad			–
Singh	– –		
Soderland		,	

Exchange Model – OEM, Object

OEM

¹⁸ Kosala, R. and H. Blockeel. Web Mining Research: A Survey (<http://facweb.cs.depaul.edu/mobasher/classes/ect584/papers/kosala.pdf>, 4.07.2011).

Goldman Widom		OEM	DataGuides
Grumbach Mecca			
Nestorov		OEM	
Toivonen		OEM	
Wang Liu		OEM	
Zaiane Han			

DataGuides²⁰:

. DataGuides

¹⁹ Kosala, R. and H. Blockeel. Web Mining Research: A Survey (<http://facweb.cs.depaul.edu/mobasher/classes/ect584/papers/kosala.pdf>, 4.07.2011).

²⁰ Goldman, R. and J. Widom. DataGuides: Enabling Query Formulation and Optimization in Semistructured Databases, (<http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.55.8594&rep=rep1&type=pdf>, 5.07.2011).

1.3.

Web Structured Mining

in-links – out-links –

(. . anchor)

()

WSM : Page Rank HITS (Hyperlink-Induced Topic Search).

Page Rank

Google.

Page Rank

N

b 1/N.

b 1/N o

Page Rank

HITS²¹, Hubs and authorities,

Page Rank

²¹ Manning, C., P. Raghavan, H. Sch tze, Hubs and Authorities (<http://nlp.stanford.edu/IR-book/html/htmledition/hubs-and-authorities-1.html>, 8.07.2011).

Page Rank HITS

1/k

k

1.4.

WUM

1.

WUM.

IP

IP

²² Srivastava, J. et. al. Web Usage Mining: Discovery and Applications of Usage Patterns from Web Data <http://portal.acm.org/citation.cfm?id=846188> (9.07.2011).

IP

²³.

e

30

²⁴.

WUM

Data Mining

WCM,

http

WUM

2.

WUM.

²³ HTTP

²⁴ Catledge, L. and J. Pitkow. Characterizing Browsing Strategies in the World-Wide Web (<http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.103.4010&rep=rep1&type=pdf>, 12.07.2011).

.
 ,
 ,
 ,
 .
 WUM.
 ,
 WUM
 ,
 ,
 .
 WUM.
 3.
 WUM.
 ,
 ,
 .
 WUM.
 ,
 ,
 .
 WUM
 .

-
 .
 -
 .
 -
 -
 -
 .
 -
 -
 -
 -
 -
 -
 -
 -
 -

II.

2.1.

”²⁵ - ” -

680 . ²⁶ . 2011 . 18,9%

• ;

• ;

• ;

• ;

• ;

• ;

• ;

• ;

• ;

• ;

²⁵ Fingar, P., . Kumar and . Sharma. Enterprise E-Commerce. Tampa: Meghan-Kiffer Press, USA, 2000.

²⁶ Rao, L., P. Morgan. Global E-Commerce Revenue To Grow By 19 Percent In 2011 To \$680B (<http://techcrunch.com/2011/01/03/j-p-morgan-global-e-commerce-revenue-to-grow-by-19-percent-in-2011-to-680b/>,10.07.2011).

Amazon.com,

•

•

•

4). WM (

, 4

WM.

?, ,,

?, ,,

?, ,,

?, ,,

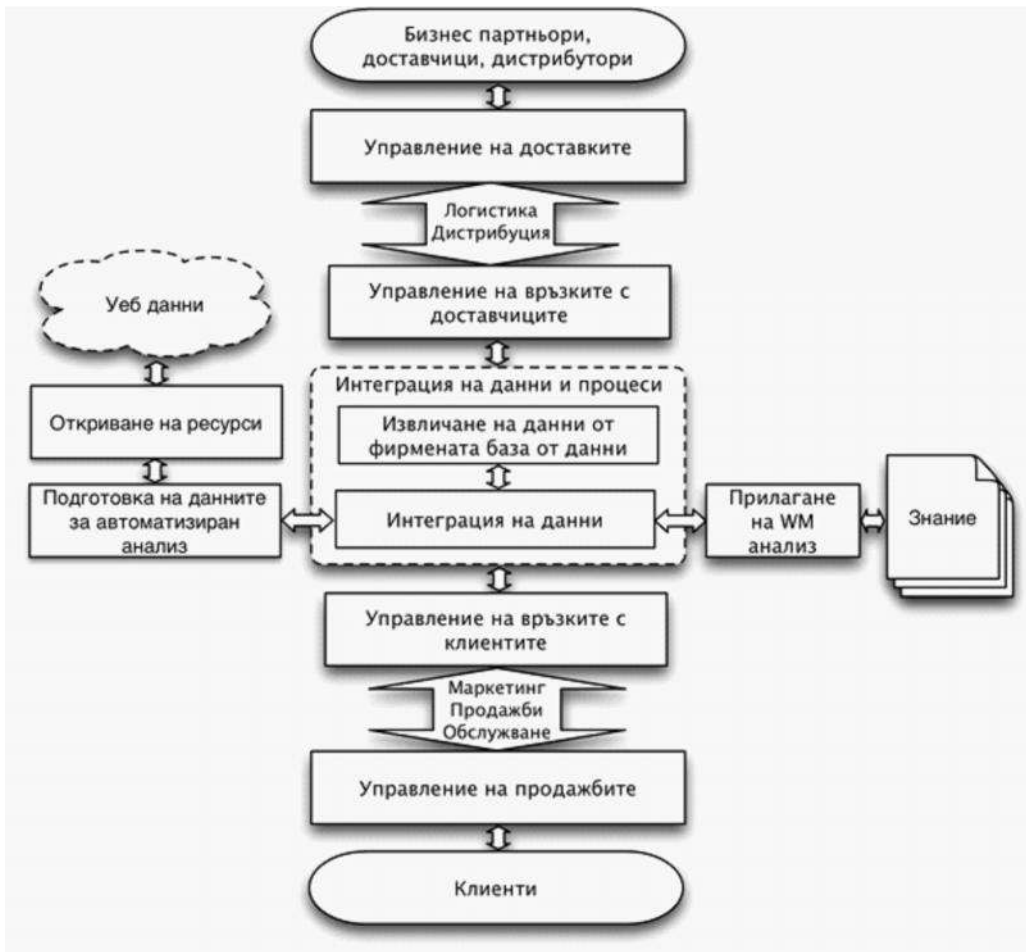
WM

WM

	Web Content Mining	Web Structure Mining	Web Usage Mining
	+	+	+
	+	-	+
	+	+	+
	+	+	+
	+	-	+
	-	+	+
	+	+	+
	+	-	+
	+	+	+
	-	+	+
-	+	+	+

2.2.

-
-
-
-



. 2.

WM а -

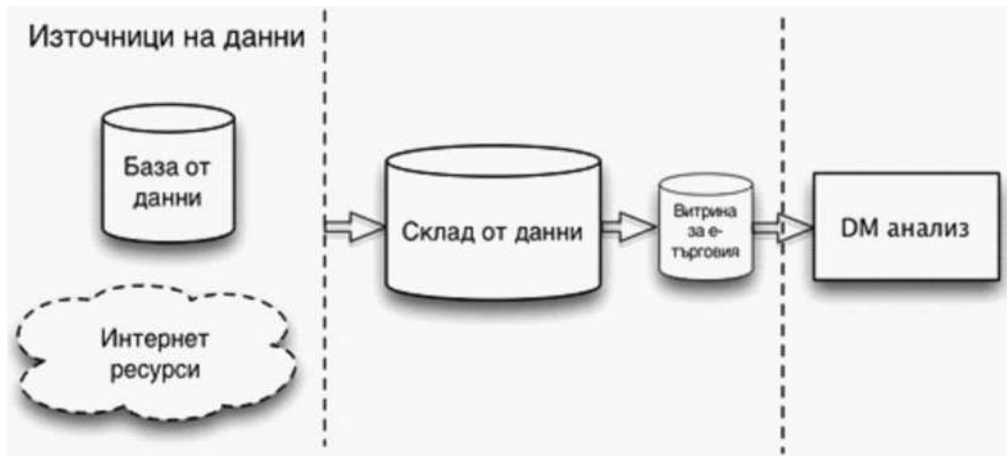
WM

mazon HP

OLAP Data Mining

(Data Marts)²⁷,

(. . . 3).
WM



. 3.

. 2

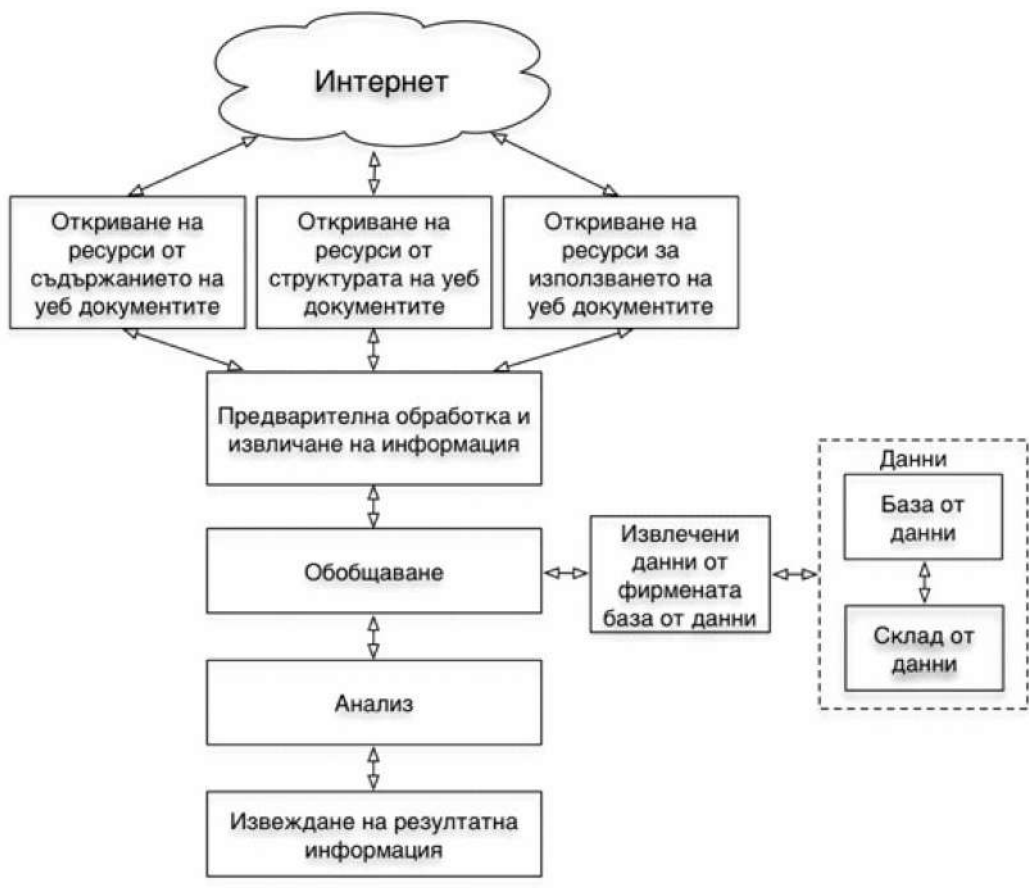
(. . 4).

²⁷



. 4.

1. , .5. :
2. () . , -
3. , , -
4. DM . -
5. , , -



. 5.

WM

2 : 6–
S 18–55 mm F 3.5–5.6 IS.

5
Canon 550D,

9– 101
Cannon EF-

Canon PowerShot.

Filter Canon UV,

Canon EOS 1100D
Canon PowerShot

1	101	6
2	101	9
3	105	12
4	204	3
...

DM WM

Cannon, -

2.3.

. 1 Data Mining
 : CRISP-DM (Cross Industry Standard Process for Data Mining),
 PMML (Predictive Model Markup Language), CWM-DM (Common Warehouse Model
 for Data Mining) PMML - 28.
 XML- Data Mining Group (DMG).

PMML-

PMML , :

-
-
-

DM

²⁸

, 2009, . 259.

Data Mining,

6 7
Data Mining -
6

Oracle Data Mining (ODM)	Oracle Database 11g. WM Oracle Text,
SAS Enterprise Miner	Java SAS Institute. log SA Link Analysis

²⁹ , 2010, .32.

³⁰ The Open Source Definition (<http://opensource.org/docs/osd>, 16.03.2011).

SPSS Modeler SPSS Text Analytics	IBM. Clementine. 2009 IBM,	SPSS Modeler SPSS Modeler. SPSS Text Analytics	SPSS
STATISTICA Data Miner			StatSoft. R

7

Carrot2	Text Mining		Java.
ELKI (Environment for DeveLoping KDD-Applications Supported by Index-Structures)		Java.	
KNIME (Konstanz Information Miner)	Eclipse,	KNIME	Java
Orange		C++ Python.	
R	Scheme. R	Ross Ihaka Robert Gentleman	
RapidMiner, YALE (Yet Another Learning Environment), RapidNet	RapidMiner	Java	RapidNet e RapidMiner.

UIMA (Unstructured Information Management Architecture)	IBM.
Weka (Waikato Environment for Knowledge Analysis)	Java Weka

-
-
-
-
-

6 7

- AeroText – text mining,
- Analog (Dr. Stephen Turner) –
- ClickTracks –
- Nihuo Web Log Analyzer –
- Surf Pattern Visual Analyzer –
- Textlyser –
- WebTrends –
- VISITaTOR – ; htminer –

– R, STATISTICA, SPSS Modeler, Rapid

Miner³¹.
 •
 •
 •
 • PMML,
 •
 • Windows, Linux Mac OS X;
 •
 DM,
 STATISTICA STATISTICA Text
 Miner STATISTICA
 Sequence, Association and Link Analysis (SAL) –
 . SPSS Modeler
 R
 Rapid Miner
 Miner
 Rapid
 SQL Server, Oracle, Access, MySQL, a CSV,
 Excel WM,

³¹ 2010 Data Miner Survey (<http://www.rexeranalytics.com/Data-Miner-Survey-Results-2010.html>, 31.08.2011); Goebel, M. and Le Gruenwald. A survey of data mining and knowledge discovery software tools (<http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.130.1456&rep=rep1&type=pdf>, 31.08.2011); Data mining. From Wikipedia, the free encyclopedia (http://en.wikipedia.org/wiki/Data_mining, 31.08.2011); Herschel, G. Magic Quadrant for Customer Data-Mining Applications, Gartner Inc. (http://www.asiaminer.com.tw/Assets/Download/Gartner_Magic_Quadrant.pdf, 7.01.2012); Haughton, D. at al. A Review of Software Packages for Data Mining. // The American Statistician, Vol. 57, 4, pp. 290–309. Most Popular Data Mining Software, (<http://www.the-data-mine.com/bin/view/Software/MostPopularDataMiningSoftware>, 31.08.2011).

STATISTICA IBM SPSS Modeler.

R Rapid Miner.

R Rapid Miner

1. , . : ,2011.
2. , . : - - ,2009.
3. , „ . , . , 2010.
4. , . : ,2011.
5. , „ . . // , . 1, 2010.
6. , . , (http://cio.bg/4380_zakonat_na_mur_e_v_sila_i_za_internet_resursite, 10.01.2012).
7. , . Data Mining // CIO, 12, 2009, .
8. - // CIO, 3, 2011 (http://cio.bg/3715_intelignentnite_biznes_analizi__tendencii_i_perspektivi, 2.07.2011).
9. BI - . // CIO, 10, 2008.
10. A. Guazzelli, M. Zeller, W. Chen, and G. Williams. PMML: An Open Standard for Sharing Models. // The R Journal, Vol. 1, 1, May 2009.
11. Ansari, S. at. al. E-Commerce and Data Mining: Architecture and Challenges (<http://ai.stanford.edu/~ronnyk/integratingEcom.pdf>, 27.09.2011).
12. Berry, M. and G. Linoff. Data Mining Techniques For Marketing, Sales, and Customer Relationship Management. Indianapolis: Wiley, 2004.
13. Buddhinath G. and D Derry. A Simple Enhancement to One Rule Classification (<http://www.buddhinath.net/OtherLinks/Documents/Improved%20OneR%20Algorithm.pdf>, 7.11.2011).
14. Catledge, L. and J. Pitkow. Characterizing Browsing Strategies in the World-Wide Web (<http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.103.4010&rep=rep1&type=pdf>, 12.07.2011).
15. Chakrabarti, S, et al. Mining the Link Structure of the World Wide Web (<http://www.cs.cornell.edu/home/kleinber/ieee99-web.pdf>, 12.01.2012).

16. Cho, Y., J. Kim. Application of Web usage mining and product taxonomy to collaborative recommendations in e-commerce. *Expert Systems with Applications* 26, 2004, p . 233–246.
17. Cooley, R, B. Mobasher and J. Srivastava. Web Mining: Information and Pattern Discovery on the World Wide Web (<http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.27.3042&rep=rep1&type=pdf>, 3.07.2011).
18. Cooley, R., P. An and J. Srivastava. Discovery of Interesting Usage Patterns from Web Data (http://www.cs.umn.edu/research/technical_reports.php?page=report&report_id=99-022, 12.01.2012).
19. Curtis, G. and Cobham, D. *Business Information Systems*. Financial Times Management, 2005.
20. Etzioni, . The World Wide Web: quagmire or gold mine? // *Communications of the ACM*, 11, 1996, p . 65–68.
21. Fayyad, U, G. Piatetsky-Shapiro and P. Smyth. From Data Mining to Knowledge Discovery in Databases (<http://www.kdnuggets.com/gpspubs/aimag-kdd-overview-1996-Fayyad.pdf>, 5.07.2011).
22. Fingar, P., . Kumar, . Sharma. *Enterprise E-Commerce*. Tampa: Meghan-Kiffer Press, USA, 2000.
23. Goebel, M., Le Gruenwald. A survey of data mining and knowledge discovery software tools (<http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.130.1456&rep=rep1&type=pdf>, 31.08.2011);
24. Goldman, R. and J. Widom. DataGuides: Enabling Query Formulation and Optimization in Semistructured Databases (<http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.55.8594&rep=rep1&type=pdf>, 5.07.2011).
25. Hand, D., . Mannila, . Smyth. *Principles of Data Mining*. MIT Press, 2001 (<http://mitpress.mit.edu/books/chapters/026208290Xchap1.pdf>, 23.08.2011).
26. Herschel, G. Magic Quadrant for Customer Data-Mining Applications, Gartner Inc. (http://www.asiaminer.com.tw/Assets/Download/Gartner_Magic_Quadrant.pdf, 7.01.2012).
27. Kadav, A., J. Kawale, P. Mitra. *Data Mining Standards* (<http://www.datamininggrid.org/wdat/works/att/standard01.content.08439.pdf>, 12.09.2011).
28. Kantardzic, . *Data Mining: Concepts, Models, Methods, and Algorithms*. John Wiley & Sons, 2003.
29. Kohavi, R. Mining E-Commerce Data: The Good, the Bad, and the Ugly, (<http://robotics.stanford.edu/~ronnyk/pakdd2001.pdf>, 12.09.2011).
30. Kosala, R. and H. Blockeel. Web Mining Research: A Survey (<http://facweb.cs.depaul.edu/mobasher/classes/ect584/papers/kosala.pdf>, 4.07.2011).
31. Larose, D. *Data mining methods and models*. Hoboken: John Wiley & Sons, 2006.
32. Manning, C., P. Raghavan, H. Sch tze. Hubs and Authorities (<http://nlp.stanford.edu/IR-book/html/htmledition/hubs-and-authorities-1.html>, 8.07.2011).
33. Markov, Z., D. Larosed. *Data Mining the Web Uncovering Patterns in Web Content, Structure, and Usage*. New Jersey: John Wiley & Sons, 2007.
34. Mikut, R., M. Reischl. *Data Mining Tools*. // *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, Vol. 1, 5, September/October 2011, p . 431–443.

35. Mladenic, D, M. Grobelnic. Feature selection for unbalanced class distribution and Naive Bayes (http://kzi.polsl.pl/~jbiesiada/prace_magisterskie/TrappLudynia/html/Bibliografia/mladenic99feature.pdf, 5.07.2011).
36. Mobasher, B., R. Cooley and J. Srivastava. Automatic Personalization Based on Web Usage Mining. // Magazine Communications of the ACM, Vol. 43, 8, 2000.
37. Rao, L., J. Morgan. Global E-Commerce Revenue To Grow By 19 Percent In 2011 To \$680B (<http://techcrunch.com/2011/01/03/j-p-morgan-global-e-commerce-revenue-to-grow-by-19-percent-in-2011-to-680b/>, 10.07.2011).
38. Shannon, Cl. Prediction and entropy of printed English. // The Bell System Technical Journal, 30 (http://www.princeton.edu/~wbialek/rome/refs/shannon_51.pdf, 12.01.2012).
39. Srivastava, J. et. al. Web Usage Mining: Discovery and Applications of Usage Patterns from Web Data (<http://portal.acm.org/citation.cfm?id=846188>, 9.07.2011).
40. Wu, X., V. Kumar. The Top Ten Algorithms in Data Mining. Chapman and Hall/CRC, 2009.
41. Zupan, B., J. Demsar. Open-Source Tools for Data Mining (<http://eprints.fri.uni-lj.si/893/1/2008-OpenSourceDataMining.pdf>, 9.09.2011).
42. 2010 Data Miner Survey (<http://www.rexeranalytics.com/Data-Miner-Survey-Results-2010.html>, 31.08.2011).
43. Data mining. From Wikipedia, the free encyclopedia (http://en.wikipedia.org/wiki/Data_mining, 31.08.2011).
44. Most Popular Data Mining Software (<http://www.the-data-mine.com/bin/view/Software/MostPopularDataMiningSoftware>, 31.08.2011).
45. RapidMiner (<http://www.rapidminer.com>, 20.11.2011).

APPLYING THE TECHNOLOGIES FOR MINING KNOWLEDGE OUT OF THE WEB IN E-COMMERCE

Assoc. Prof. Dr Snezhana Salova

Abstract

Of great significance for business is knowledge that can be derived from unstructured information contained in Internet sources, such as text, hyperlinks, tags, log files, etc.

The paper presents a study of the process of extraction of useful knowledge from web resources. There is proposed a model of a system of e-commerce, in which there are integrated the activities of mining knowledge out of Internet sources. There has been developed a functional matrix of application of the various kinds of web mining for the sphere of e-commerce. There are also defined the required software tools for the realization of the presented model.

ANWENDUNG DER TECHNIKEN ZUR INFORMATIONSGEWINNUNG IM INTERNET IM ELEKTRONISCHEN HANDEL

Doz. Dr. Snezhana Salova

Zusammenfassung

Das Wissen, das die Subjekte der Wirtschaft aus den nicht-strukturierten Informationen im Internet gewinnen können, ist von wesentlicher Bedeutung für die Wirtschaft. Es handelt sich um Texte, Hyperlinks, Tags, Log-Dateien u.a.

In der Studie wird der Vorgang der Gewinnung nützlichen Wissens aus Web-Ressourcen untersucht. Ein Model eines Systems elektronischen Handels wird vorgeschlagen, in dem die Tätigkeiten zur Gewinnung von Wissen aus dem Internet mit dem Handel selbst integriert sind. Eine funktionelle Matrix zum Einsatz verschiedener Arten von Web Mining für den Bereich des elektronischen Handels wird dargestellt. Software-Mittel zur Verwirklichung des dargestellten Models werden definiert.

. -

Web Mining

.....	198
I.	199
1.1.	199
1.2.	206
1.3.	211
1.4.	212

II.	215
2.1.	215
2.2.	-	220
2.3.	225
	230
	231
	233
	234
	234