

004.912

« »,

WEB-

Web Content Mining,

, *Web Content Mining,*

HTML.

Web Content Mining –

Web Content Mining,

Web Mining,

Web Structure Mining Web Usage Mining.

Web Content Mining,

Mining

Web Content

[1].

Web-

HTML-

[2].

Natural Language

Processing (NLP),

[3])

1.

<concept >[.,] , <subconcept > s+
<concept >[.,] , <subconcept > s+

() ,

(NN)

(NP),

<concept>, <concept_r> <supconcept>,

html-
Web Content

Mining :

NLP,

<subconcept >

<NN> <NP>

<concept>, NP <concept>

<NN> <NP>,

<subconcept > :

< concept >[.,] , < subconcept > s +

< concept > s * such as < subconcept >

: “... ”

“ ” “ ”

“ ” “ ”

“ ”

< concept _r > s * [.,]e.g., (< concept > s+)+,

< concept _r > [.,] [,](< concept > s+)+

<concept_r>

<NN>

<NP>,

<concept>

“There are many Data Mining techniques, e.g. Clustering, Classification, Data Warehouses, Web Mining etc.” “Data Mining techniques”,

{“Clustering”, “Classification”, “Data Warehouses”, “Web Mining”}.

1. Morinaga K. Mining product reputations on the Web / K.Morinaga, K. Yamanishi , Tateishi and T. Fukushima // Procs. of KDD. – 2002. – P. 341–349.

2. Wang Y. Web Mining and Knowledge Discovery of Usage Patterns / Y. Wang–2000. <https://cs.uwaterloo.ca/~tozsu/courses/cs748t/surveys/wang.pdf>.

3. Broder A. Robust classification of rare queries using web knowledge / Broder A., Fontoura M., Gabrilovich E. et al. // Proceedings of the 30th Annual International ACM SIGIR Conference on Research and development in information retrieval (SIGIR’07), 2007. – P. 231–238.

4. Lan Yi. Eliminating Noisy Information in Web Pages for Data Mining / Yi Lan, Liu Bing, Li Xiaoli // Proceedings of the ninth ACM SIGKDD International Conference on Knowledge Discovery & Data Mining (KDD-2003), Washington, 2003. – P. 296 – 305.

5. // , 2012. – . 38. – . 239 – 245.

12.11.2014

WEB-*Mining,**Web Content**, Web Content Mining,***FEATURES EXTRACTION AND IDENTIFICATION OF KNOWLEDGE WEB-CONTENT**

N.F. Hayrova, Ajit Pratap Singh Gautam

In the article the features of knowledge mining and knowledge identification of web-pages have been considered. The new kind technology of Web Content Mining has been elaborated. The technology is based on the method of extraction of semantic concepts from textual information and includes the steps: exarticulation of the main page-content, extraction of the semantic concepts and the content analysis. At the stage of content analysis regular expressions have been used. The regular expressions allow to manifestly distinguish relationships of the representation and taxonomy between concepts of the webpage. As elements of regular expressions were used nouns, nouns groups and special lexical constructs.

Keywords: *the identification of knowledge, Web Content Mining, regular expressions, taxonomy, attitude of representativeness.*