

**MINISTERUL EDUCAȚIEI, CULTURII ȘI CERCETĂRII AL
REPUBLICII MOLDOVA**

**Universitatea Tehnică a Moldovei
Facultatea Calculatoare, Informatică și Microelectronică
Departamentul Ingineria Software și Automatică**

Admis la susținere
Șef departament:
Ion Fiodorov, conferențiar universitar, doctor în informatică

“ ” _____ 2021

**Instrument de analiză a datelor textuale utilizând tehnici
Data Mining**

Teză de master

Student: Scutelnic Bogdan, TIA-191M

Conducător: Rogovschi Nicoleta, conf. univ.dr.

Chișinău, 2021

ADNOTARE
la teză pentru obținerea titlului de master cu tema
La teză de master “INSTRUMENT DE ANALIZĂ A DATELOR TEXTULAE UTILIZÂND
TEHNICI DATA MINING”

Scutelnic Bogdan
Specializarea Tehnologia Informației pentru Afaceri

Actualitatea temei. Data Mining este un proces de a găsi modele potențial utile din seturi de date imense. Este o abilitate multidisciplinară care folosește învățarea automată, statistici și AI pentru a extrage informații pentru a evalua probabilitatea evenimentelor viitoare. Statisticile derivate din Data Mining sunt utilizate pentru marketing, detectarea fraudei, descoperirea științifică etc.

O dată cu dezvoltarea rapidă a tehnologiilor informaționale a apărut necesitatea dezvoltării sistemelor informaționale, pentru a programa sisteme soft care ar putea stoca cantități enorme de informații și ar putea reda prin comenzi simple și ușor de asimilat date exacte despre conținutul bazei de date, precum și prelucrarea unor altor tipuri de date nefiind obligatoriu să prelucreze date sub formă de tabele dar pot fi date grafice sau date de avertizare sau de alt tip din domeniul informaticii, care sigur este domeniul cel mai important la ziua de azi precum și cel mai cercetat care datorită posibilități vaste pe care le poate îndeplini un soft combinat cu hard. Astfel viitorul zilei de mâne al oricărui domeniu nu poate fi conceput fără folosirea unui soft al domeniului respectiv.

Structura lucrării: lucrarea este alcătuită din introducere, trei capitole, concluzii și bibliografie.

În capitolul unu sunt analizate principalele noțiuni despre data mining, text minig, depozite de date, precum și noțiunea de Big Data.

În capitolul doi sunt prezentate tehnologiile și metodele pentru explorarea datelor prin tehnici specifice proceselor data mining.

În capitolul trei este prezentată aplicația care extrage date utilizând tehnicile data minig.

Cuvinte-cheie: Data mining, text mining, depozit de date, descoperire de resurse, statistica datelor.

Scopul lucrării: Elaborarea unui instrument care va permite analiza datelor text stocate în baza prelucrării statistice și a tehnicilor Data Mining.

Tehnologii utilizate: data mining, statistica datelor, limbajul R+.

În lucrare sunt studiate metode de prelucrare a textului și crearea modulelor de determinare a modelelor de detectare a conținutului. Este un lucru enorm în cadrul unui sistem de căutare centralizat sau pentru un sistem inteligent. A fost creat un instrument de anilza a datelor textuale cu ajutorul tehnicilor Data Mining si desigur intrumentul a fost verificat, fiind prezentate rezultate unei anilize.

ANNOTATION

to the Master's thesis with the theme

For my master degree thesis with the topic " TEXTUAL DATA ANALYSIS TOOL USING DATA MINING TECHNIQUES"

Scutelnic Bogdan

Specialization Business Information Technology

The actuality of the subject. Data Mining is a process of finding potentially useful models from huge datasets. It is a multidisciplinary skill that uses machine learning, statistics and AI to extract information to assess the likelihood of future events. Statistics derived from Data Mining are used for marketing, fraud detection, scientific discovery, etc.

With the rapid development of information technologies came the need to develop information systems to program software systems that could store huge amounts of information and could render through simple and easy to assimilate accurate data about the contents of the database, as well as processing other types of data are not required to process data in the form of tables but can be graphical data or warning or other data in the field of informatics, which is certainly the most important field today and the most researched due to the possibilities vast that a software combined with hardware can accomplish. Thus, the future of tomorrow of any field cannot be conceived without the use of software of that field.

Structure of the thesis: the paper consists of introduction, three chapters, conclusions and bibliography.

Chapter one analyzes the main notions about data mining, text mining, data warehouses, as well as the notion of Big Data.

Chapter two presents the technologies and methods for data exploration through techniques specific to data mining processes.

Chapter three presents the application that extracts data using mining data techniques.

Keywords: Data mining, text mining, data storage, resource discovery, data statistics.

Purpose of the paper: Development of a tool that will allow the analysis of text data stored based on statistical processing and Data Mining techniques.

Technologies used: data mining, data statistics, R + language.

The paper studies methods of word processing and the creation of modules for determining content detection models. It's a huge thing in a centralized search system or for a smart system. A text data analysis tool was created using Data Mining techniques and of course the tool was verified, presenting the results of an analysis.

CUPRINS

INTRODUCERE	8
1. ANALIZA DOMENIULUI DE CERCETARE	10
1.1 Procesul de data mining	11
1.2 Text mining.....	16
1.3 Depozite de date.....	20
1.4 Conceptul de BigData.....	22
2. TEHNOLOGII INFORMATICE INTELIGENTE DE ACCESARE	24
MULTIDIMENSIONALĂ A BAZELOR ȘI DEPOZITELOR DE DATE	24
2.1 Tehnologia OLAP.....	25
2.2 Modele de descoperire de resurse prin extragerea datelor.....	28
2.3 Tehnici de extragere a textului.....	32
3. REALIZAREA APLICAȚIEI	36
3.1 Alegerea strategie de proiectare	37
3.2 Alegerea instrumentarului	42
3.3 Descrierea algoritmului.....	47
3.4 Descrierea aplicației.....	54
CONCLUZII	63
REFERINȚE BIBLIOGRAFICE	64

INTRODUCERE

Domeniul gestiunii informațiilor abordează problema organizării, stocării și regăsirii în timp util a datelor de care dispunem despre un anumit subiect. Neîndoind că fiecare dintre noi se confruntă continuu cu această problemă. Agenda personală, cartea de telefoane, sunt mijloace simple de organizare, păstrare și regăsire a datelor de care avem nevoie la un moment dat.

"Avem foarte multe date colectate, ce să facem cu ele acum?" Această problemă a devenit una obișnuită în multe organizații mari. Informația digitală este ieftin de obținut și relativ ieftin de stocat. Dar care este scopul stocării unei cantități de date atât de mari? În afara argumentului legat de modalitățile convenabile de stocare a datelor în format electronic, mai există și altul: firmele colectează date deoarece managerii "simt" că aceste date reprezintă un activ valoros, care ar putea fi folosit la un moment dat. În institutele de cercetare, datele se referă la observații asupra fenomenelor aflate sub studiu, culese cu atenție. În organizațiile economice, datele și informațiile se referă la piețe, concurenți, furnizori, distribuitori, clienți, precum și la informații interne legate de procesele de producție, organizarea muncii, depozitare etc.

Pentru multe firme din Occident, procesul de explorare și analiză a datelor nu este unul nou. Dar programele de *data mining* permit realizarea acestor analize mult mai rapid și, potențial, cu o eficiență sporită. În multe cazuri, tehnicile de explorare a datelor permit ca datele colectate pentru un anumit scop să poată fi folosite în multe alte scopuri.

Ca urmare a perfecționării modalităților de colectare și stocare a datelor, foarte multe companii mari s-au trezit deodată că stau pe un munte de date cărora ar trebui să le găsească o utilizare. De exemplu, companiile furnizoare de cărți de credite înregistrează permanent datele referitoare la tranzacțiile comerciale, supermarketurile înregistrează date privind cumpărăturile, utilizarea cupoanelor de reduceri ș.a.m.d. Apariția necesității de a valorifica toate aceste date era doar o chestiune de timp.

Chiar dacă proiectanții interfețelor și a conținutului web-ului includ comportamentul utilizator și câteva deprinderi în unele motoare de căutare, utilizatorii pot avea dificultăți în efectuarea de căutări pe web deoarece aceștia nu sunt capabili să formuleze interogările corecte care să reducă numărul de rezultate obținute și să crească calitatea acestora. De obicei interfața utilizator este dificil de utilizat, funcțiile de organizare ierarhică sunt aplicate static și informațiile de pe web sunt foarte rar structurate și organizate pentru o recunoaștere rapidă.

Pentru a rezolva problemele prezentate mai sus acest domeniu de cercetare a crescut continuu în direcții cum ar fi algoritmi, strategiile și arhitecturile. Sute de idei au fost testate, unele fiind implementate iar altele existând doar ca și prototipuri. Unele dintre aceste idei sunt: organizarea documentelor în categorii predefinite, îmbunătățirea modului de prezentare a rezultatelor, monitorizarea unei pagini specifice, ajutarea utilizatorului în formularea interogării corecte și dezvoltarea de programe pentru filtrarea rezultatelor căutării.

O dată cu dezvoltarea rapidă a tehnologiilor informaționale a apărut necesitatea dezvoltării sistemelor informaționale, pentru a programa sisteme soft care ar putea stoca cantități enorme de informații și ar putea reda prin comenzi simple și ușor de asimilat date exacte despre conținutul bazei de date, precum și prelucrarea unor altor tipuri de date nefiind obligatoriu să prelucreze date sub formă de tabele dar pot fi date grafice sau date de avertizare sau de alt tip din domeniul informaticii, care sigur este domeniul cel mai important la ziua de azi precum și cel mai cercetat care datorită posibilități vaste pe care le poate îndeplini un soft combinat cu hard. Astfel viitorul zilei de mâine al oricărui domeniu nu poate fi conceput fără folosirea unui soft al domeniului respectiv.

REFERINȚE BIBLIOGRAFICE

- [1] Andone I. Tugui, A., s.a. Dezvoltarea sistemelor inteligente în economie, Ed. Economică, București, 2011
- [2] Cipolla, E. T. Data Mining: Techniques to Gain Insight Into Your Data, Enterprise Systems Journal, Decembrie 2005;
- [3]
- [4] Guandong Xu, Zanchun Yhang, Lin Li, Web Mining and Social Networking Techniques and Applications, USA:Springer, 2011
- [5] J. Palau, M. Montaner, B. Lopez and J.L. de la Rosa. Collaboration analysis n the recommender system using social networks. In CIA, pages 137-151,2004
- [6] J.M. Kleinberg. Authoritative sources in a hyperlinked environment. In Proc. Of the Ninth Annual ACM-SIAM Symposium on Discrete Algorithms (SODA '98), pages 668-677, 1998
- [7] Zaiane O., Han J.: WebML: Querying the World Wide Web for resources and knowledge. In: Workshop on Web Information and Data Management WIDM98, Bethesda, 1998, 9-12.
- [8] Yang, Q., Zhang, H.H. and Li, I.T, 2001. Mining Web logs for prediction models in www caching and prefetching, Seventh ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, San Francisco, CA, USA, August 26-29, pp. 473-478.
- [9] Yan, T.W., Jacobsen, M. Garcia-Molina, H. and Dayal, U., 1996, Knowledge Discovery from users web page navigation, Seventh International Workshop on Research Issues in Data Engineering, Birmingham, England, aprilie 7-8, pp 20-29.
- [10] Berry, M., Linoff, G.: Data Mining Techniques for Marketing, Sales and Customer Support, John Wiley and Sons, Chichester (1997)
- [11] Dunham, M.H.: Data Mining : Introductory and Advanced Topics. Prentice Hall, Pearson Education Inc. (2003)
- [12] Prinzie, A. , Van den Poel, D.: Investigating Purchasing Patterns for Financial Services using Markov, MTD and MTDg Models. In: Working Papers of Faculty of Economics and Business Administration, Ghent University, Belgium 03/213 (2013)
- [13] Agrawal, R., Srikant, R.: Mining sequential patterns, *International Conference on Data Engineering(ICDE'95)*, Taipei, Taiwan, pp. 3-14, martie 1995
- [14] Inmon, W.H. Buiding the Data Warehouse, John Wiley & Sons, New York, 2005
- [15] Turban, E., Jay Aronson J. Decision Support Systems and Intelligent Systems, Ed. Prentince Hall, New Jersey, USA, 2011