

DATA MINING
1. Date despre unitatea de curs

Facultatea	Calculatoare, Informatică și Microelectronică				
Catedra/Departamentul	Ingineria Software și Automatică				
Ciclul de studii	Studii superioare de master, ciclul II				
Programul de studii	Tehnologia informației pentru afaceri				
Anul de studii	Semestrul	Tip de evaluare	Categoria formativă	Categoria de opționalitate	Credite ECTS
II (învățământ cu frecvență);	3	E-examen	S – unitate de curs de specialitate	A - unitate de curs opțională	5

2. Timpul total estimat

Total ore în planul de învățământ	Din care				
	Ore auditoriale		Lucrul individual		
	Curs	practice	Proiect de an	Studiul materialului teoretic	Pregătire aplicații
150	20	20	-	110	-

3. Precondiții de acces la unitatea de curs

Conform planului de învățământ	Analiza și exploratorie a datelor, managementul datelor, fundamente ale tehnologiei informației
Conform competențelor	Posedarea cunoștințelor medii de utilizare a instrumentelor Microsoft Office. Posedarea abilităților de folosire a calculatorului. Competențe de înțelegere și aplicare a formulelor matematice.

4. Condiții de desfășurare a procesului educațional pentru

Curs	Pentru prezentarea materialului teoretic în sala de curs este nevoie de proiector, PC/laptop și acces la internet. Nu vor fi tolerate întârzierile studenților, precum și convorbirile telefonice în timpul cursului
Laborator/Seminar	Studenții vor perfectă rapoarte conform condițiilor impuse de indicațiile metodice. Termenul de predare a lucrării de laborator - o săptămână după finalizarea acesteia. Pentru predarea cu întârziere a lucrării aceasta se depunceață cu 1pct./săptămână de întârziere.

5. Competențe specifice acumulate

Competențe profesionale	C1 Operarea cu concepte și metode ale domeniului Tehnologiei Informației C2 Aspecte organizaționale și informaționale ale sistemelor C3 Modelarea sistemelor informaționale complexe și implementarea lor prin sisteme informatice C5 Managementul produselor și al serviciilor TIC în concordanță cu cerințele pieței
Competențe transversale	CT1. Aplicarea principiilor, normelor și valorilor eticii profesionale CT2. Identificarea, descrierea și derularea activităților organizate într-o echipă cu dezvoltarea capacităților de comunicare și colaborare, dar și cu asumarea diferitelor roluri (de execuție și conducere) CT3. Demonstrarea spiritului de inițiativă și acțiune pentru actualizarea propriilor cunoștințe profesionale, economice și de cultura organizațională

6. Obiectivele unității de curs

Obiectivul general	Utilizarea metodelor de căutare și extragere a informației în domeniul dobândirii datelor (data mining) permite reducerea timpului de rezolvare a problemelor informaționale, mărește precizia dobândirii informației nestructurate și oferă soluții practic admisibile acolo unde soluții exacte sînt imposibil de abordat. Metode și tehnici de text mining trebuie să de vină unele din instrumentele curente cercetări
Obiectivele specifice	La terminarea cursului „Sisteme de cautare a informației” studentul trebuie să cunoască: * Tehnologiile construirii șabloanelor (modelelor „patterns”) pentru formarea diferitor structuri de text. * Metode statistice și deterministice pentru sumarizarea textului. * Aplicarea metodelor de căutare a informației (Information Retrieval) pentru rezolvarea problemelor „Text Mining” * Tehnologii multivariante (word corelation, word covariance și component analysis) pentru procesarea textului. * Aplicarea metodelor de clasterizare a textului pentru dobîndirea datelor nestructurate. * Utilizarea limbajului PERL pentru rezolvarea problemelor „Text Mining”

7. Conținutul unității de curs

Tematica cursului	Numărul de ore învățămînt cu frecvență
T.1. Text patterns. Introduction. Regular Expression. Finding words in a text. Decomposing Poe’s „The Tell- Tale Heart” into words.	1
T.2. Text Paterns. Simple concordance First Attempt at Extracting Sentences. Regular expression Odds and Ends	1
T.3. Quantitative Text Summaries . Introduction Scalars, Interpolation and conect in Perl. Word Lengths in Poe’s „The Tell- Tale Heart”. Array and Functions.	1
T.4. Quantitative Text Summaries. Using Hash. Zipf’ law for „A Christmas Carol”. Word Anagrams. Finding Words in setof letters. Complex Data Structures	1
T.5. Probability and text Sampling. Probability. Conditional probability. Mean and Variance of random, Variable. The Bag-of-Words Model for Poe’s „The Blaak Cat	1
T.6. Applying Information. Retrieval to text mining. Introduction. Counting letters and Words. Text counts and vectors. The Term-Document Matrix Applied to Poe.	1
T.7. Applying Information. Retrieval to text mining. Matrix Multiplication. Function of Counts. Document Similarity. Invers Document Frequency. Poe Story Angles Revisired	1
T.8. Concordance Lines and Corpus Linguistics. Text Sampling. Statistical survey sampling. Corpus as Baseline. Sorting concordance lines.	1
T.9. Concordance Lines and Corpus Linguistics. Aplication: Word Usage Differentces between London and Sheley. Word morphology of adverbs. Collocations and Concordance Lines. More ways to sort concordance lines.	2
T.10. Multivariance Techniques with text. Basic Statistics. Z-scores appllied to poe. Correlation and Cosines. Correlation and Covariances. 2 by 2 Correlation and Covariance	2
T.11. Multivariate Techniques with text. Principal Components Analysis (PCA) Finding the principal Components. PCA Applied to the N Poe short stories	2
T.12. Text Clustering. Two variable Example of K-Means. K-Means with language R. Poe Cluster using Eight Prononns.	2
T.13. Text Clustering. Clustering Poe using principal components. Hierarchical clustering of Poe’s short Stories. Decision Trees and Overfinding.	2
T.14. A Simple of additional Topics. Perl Module: module for number words; the stop-moduler the senteces segmentation module; an object oriented Module for tagging. Miscellaneous modules.	2
Conclusions	
Total ore curs:	20

Tematica lucrărilor de laborator/seminarelor	Numărul de ore
	învățământ cu frecvență
L1. Aplicarea instrucțiunilor de bază a limbajului Perl pentru procesarea propozițiilor (array, lengths, function: split, join, substr, pos ș.a.	4
L2. Construirea mini-mașinii de cautare a informației în baza textelor electronice date	4
L3. Testarea și obținerea rezultatelor funcționării programelor în limbajul Perl pentru concordance sofisticate în baza modelelor specifice pentru procesarea structurilor lingvistice date.	4
L4. Elaborarea algoritmului clasterizării textelor în baza metodei K-Means	4
L5. Oformarea și prezentarea finală a rapoartelor pentru lucrările de laborator 1-6	4
Total ore practice:	20

8. Referințe bibliografice

Principale	1. Ronen Feldman. The Text Mining Handbook. Cambridge. 2007. 2. Roger Bilisoly . Practical Text Mining With Perl. Canada 2008. 3. Chris Manning, Hinrich Schütze. 1999. Foundations of Statistical Natural Language Processing. MIT Press, Cambridge, MA. 2. Christopher D. Manning, Prabhakar Raghavan and Hinrich Schütze, Introduction to Information Retrieval, Cambridge University Press. 2008 4. Stuart Russel, Peter Norving , Artificial Intelligence a Modern Approach. New Delhi.2005
Suplimentare	1. Roger Bilisoly. Practical text mining with Perl. Willey. 2008. 295p. 2. Cristopher Dmeaning Prabhacar, Ragharan. Hinrich Schutze. An introduction to Information Retrieval. Cambridge University Press. 2009. 3. Randal L. Schwartz & Tom Christianseu. Learning Perl. O'REILLY. 1997. 267p.

9. Evaluare

Periodică		Curentă	Studiu individual	Proiect/teză	Examen
EP 1	EP 2				
10%	10%	10%	30%	-	40%
Standard minim de performanță Prezența și activitățile la prelegeri și lucrări practice. Obținerea notei minime de „5” la fiecare dintre lucrări și examen					

10. Criterii de evaluare

Activitate	Componente evaluare	Metodă de evaluare, Criterii de evaluare	Pondere în nota finală a activității	Ponderea în evaluarea disciplinei
Evaluare periodică I	Conținut teoretic, teme 1-7	Test pe MOODLE	100%	10%
Evaluare periodică II	Conținut teoretic, teme 8-14	Test pe MOODLE	100%	10%
Evaluare curentă	Lucrări practice	Discuții în cadrul lecțiilor practice	50%	10%
		Dosar completat cu Rapoarte pentru fiecare Studiu de caz în discuție	50%	
Studiul individual	Teme individuale	Prezentare/discurs public	100%	30%
Evaluarea finală	Conținut teoretic și practic	Test pe MOODLE. Notare conform baremului	100%	40%