

**PRASAD V. POTLURI SIDDHARTHA INSTITUTE OF TECHNOLOGY**

(Autonomous)  
Kanuru, Vijayawada-520007

**DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING (AI&ML)****III B Tech – I Semester****Information Retrieval Lab**

<b>Course Code</b>	23AM3551	<b>Year</b>	III	<b>Semester</b>	I
<b>Course Category</b>	PC	<b>Branch</b>	CSE (AI&ML)	<b>Course Type</b>	Practical
<b>Credits</b>	1.5	<b>L-T-P</b>	0-0-3	<b>Prerequisites</b>	Python Programming
<b>Continuous Internal Evaluation</b>	30	<b>Semester End Evaluation</b>	70	<b>Total Marks</b>	100

**Course Outcomes****Upon Successful completion of course, the student will be able to**

<b>CO1</b>	Understand basic techniques in information retrieval and text preprocessing including document representation, stopword removal, and stemming.	<b>L2</b>
<b>CO2</b>	Apply core algorithms for text classification, clustering, and web crawling using standard datasets and real-world content using tools	<b>L3</b>
<b>CO3</b>	Analyze large-scale document collections and web content using indexing, graph-based ranking, and topic modeling techniques.	<b>L4</b>
<b>CO4</b>	Evaluate advanced techniques such as LSI, topic-specific PageRank, and Twitter mining to extract meaningful insights and trends.	<b>L5</b>

**Contribution of Course Outcomes towards achievement of Program Outcome & Strength of correlation (3: High, 2: Medium, 1: Low)**

	PO1	PO2	PO3	PO4	PO5	PO6	PO7	PO8	PO9	PO10	PO11	PSO1	PSO2
<b>CO1</b>	2								2				
<b>CO2</b>	3	3			3							3	
<b>CO3</b>	3	3										3	
<b>CO4</b>		3									2		

**PRASAD V. POTLURI SIDDHARTHA INSTITUTE OF TECHNOLOGY**

(Autonomous)  
Kanuru, Vijayawada-520007

**DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING (AI&ML)****III B Tech – I Semester**

<b>Syllabus</b>		
<b>Exp. No.</b>	<b>Contents</b>	<b>Mapped CO</b>
<b>1</b>	Representation of a Text Document in Vector Space Model and Computing Similarity between two documents.	CO1 to CO4
<b>2</b>	Pre-processing of a Text Document: stop word removal and stemming	CO1 to CO4
<b>3</b>	Construction of an Inverted Index for a given document collection comprising of at least 50 documents with a total vocabulary size of at least 1000 words.	CO1 to CO4
<b>4</b>	Classification of a set of Text Documents into known classes (You may use any of the Classification algorithms like Naive Bayes, Max Entropy, Rochio's, Support Vector Machine). Standard Datasets will have to be used to show the results.	CO1 to CO4
<b>5</b>	Text Document Clustering using K-means. Demonstrate with a standard dataset and compute performance measures- Purity, Precision, Recall and F-measure.	CO1 to CO4
<b>6</b>	Crawling/ Searching the Web to collect news stories on a specific topic (based on user input). The program should have an option to limit the crawling to certain selected websites only.	CO1 to CO4
<b>7</b>	To parse XML text, generate Web graph and compute topic specific page rank	CO1 to CO4
<b>8</b>	Implement Matrix Decomposition and LSI for a standard dataset.	CO1 to CO4
<b>9</b>	Mining Twitter to identify tweets for a specific period (and/or from a geographical location) and identify trends and named entities.	CO1 to CO4
<b>10</b>	Implementation of PageRank on Scholarly Citation Network	CO1 to CO4

**Learning Resources****Text Books**

# PRASAD V. POTLURI SIDDHARTHA INSTITUTE OF TECHNOLOGY

(Autonomous)  
Kanuru, Vijayawada-520007

## DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING (AI&ML)

### III B Tech – I Semester

- |   |
|---|
| <ol style="list-style-type: none"> <li>1. <b>Modern Information Retrieval</b>, Ricardo Baeza-Yates &amp; Berthier Ribeiro-Neto, 2nd ed., 1999, Addison-Wesley</li> <li>2. <b>Information Storage and Retrieval Systems: Theory and Implementation</b>, Gerald J. Kowalski &amp; Mark T. Maybury, 2nd ed., September 30, 2000, Springer (Kluwer Academic)</li> </ol> |
|---|

#### References

- |   |
|---|
| <ol style="list-style-type: none"> <li>1. <b>Information Retrieval: Implementing and Evaluating Search Engines</b>, Stefan Büttcher, Charles L. A. Clarke, Gordon V. Cormack, 1st Edition, 2010, The MIT Press.</li> <li>2. <b>Information Retrieval: A Health and Biomedical Perspective</b>, William Hersh, 4th Edition, 2022, Springer.</li> </ol> |
|---|

#### E-Recourses and other Digital Material

- |  |
|--|
| <ol style="list-style-type: none"> <li>1. <a href="https://www.analyticsvidhya.com/blog/2015/04/information-retrieval-system-explained/">https://www.analyticsvidhya.com/blog/2015/04/information-retrieval-system-explained/</a></li> <li>2. <a href="https://www.analyticsvidhya.com/blog/2017/11/information-retrieval-using-kdtree/">https://www.analyticsvidhya.com/blog/2017/11/information-retrieval-using-kdtree/</a></li> <li>3. <a href="https://medium.com/analytics-vidhya/information-retrieval-part-1-extracting-webpages-a9d0b715535d">https://medium.com/analytics-vidhya/information-retrieval-part-1-extracting-webpages-a9d0b715535d</a></li> </ol> |
|--|