



**BITS Pilani**  
Pilani | Dubai | Goa | Hyderabad

## **BIRLA INSTITUTE OF TECHNOLOGY & SCIENCE, PILANI WORK INTEGRATED LEARNING PROGRAMMES**

### **COURSE HANDOUT**

#### **Part A: Content Design**

<b>Course Title</b>	Information Retrieval
<b>Course No(s)</b>	SS ZG537
<b>Credit Units</b>	4 (Unit split between Class Hours + Lab/Design/Fieldwork + Student preparation) Ex. 1-1-2, (total 4 units or credits) ie 1 unit for class room hours, 1 unit for lab hours, 2 units for student preparation. Typically 1 unit translates to 32 hours
<b>Course Author</b>	Poonam Goyal
<b>Version No</b>	1.0
<b>Date</b>	

#### **Course Objectives**

<b>No</b>	<b>Course Objective</b>
<b>CO1</b>	To understand structure and organization of various components of an IR system
<b>CO2</b>	To understand information representation models, term scoring mechanisms, etc. in the complete search system
<b>CO3</b>	To understand architecture of search engines, crawlers and the web search
<b>CO4</b>	To understand cross lingual retrieval and multimedia information retrieval

#### **Text Book(s)**

<b>T1</b>	C. D. Manning, P. Raghavan and H. Schutze. Introduction to Information Retrieval, Cambridge University Press, 2008. <a href="http://nlp.stanford.edu/IR-book/">http://nlp.stanford.edu/IR-book/</a>
-----------	---

#### **Reference Book(s) & other resources**

<b>R1</b>	Modern Information Retrieval, Ricardo Baeza-Yates and Berthier Ribeiro-Neto, Addison-Wesley, 2000. <a href="http://people.ischool.berkeley.edu/~hearst/irbook/">http://people.ischool.berkeley.edu/~hearst/irbook/</a>
<b>R2</b>	Ricci, F.; Rokach, L.; Shapira, B.; Kantor, P.B. (Eds.), Recommender Systems Handbook. 1st Edition., 2011, 845 p. 20 illus., Hardcover, ISBN: 978-0-387-85819-7
<b>R3</b>	Cross-Language Information Retrieval by Jian-Yun Nie Morgan & Claypool Publisher series 2010
<b>R4</b>	Multimedia Information Retrieval by Stefan M. Rüger Morgan & Claypool Publisher series 2010.

R5	Information Retrieval: Implementing and Evaluating Search Engines by S. Buttcher, C. Clarke and G. Cormack, MIT Press, 2010.
R6	Web Data Mining: Exploring Hyperlinks, Contents, and Usage Data by B. Liu, Springer, Second Edition, 2011.

### **Modular Content Structure**

1. Introduction
  - 1.1. Information Retrieval
  - 1.2. Basic Search Model
2. Basic Information Retrieval Concepts
  - 2.1. Boolean Retrieval
  - 2.2. Dictionaries and Tolerant Retrieval
  - 2.3. Index Construction and Compression
3. Vector Space Model
  - 3.1. Scoring, Term Weighting
  - 3.2. The Vector Space Model for Scoring
4. Text Mining
  - 4.1. Text Classification
  - 4.2. Vector Space Classification
  - 4.3. Text Clustering
5. Web Search
  - 5.1. Web Search Basics
  - 5.2. Web Crawlers and Indexes
  - 5.3. Link Analysis
6. Cross Lingual Retrieval
  - 6.1. Language Problems in IR
  - 6.2. Approaches for CLIR
7. Multimedia Information Retrieval
  - 7.1. Multimedia Search Technologies
  - 7.2. Content Based Retrieval
8. Recommender Systems
  - 8.1. Collaborative and Content Based RS

**Learning Outcomes:**

No	Learning Outcomes
LO1	Students will gain understanding about an information retrieval system as a whole and about its components.
LO2	Students will have knowledge about the design issues and their solutions of different type of models including Boolean, vector space etc.
LO3	Students will have detailed understanding about text indexing, mining, weighting schemes etc.
LO4	Students will acquire knowledge about cross lingual and multimedia information retrieval.
LO5	With the acquired knowledge students will be able to design and build different kind of information retrieval systems.

**Part B: Contact Session Plan**

<b>Academic Term</b>	
<b>Course Title</b>	Information Retrieval
<b>Course No</b>	
<b>Lead Instructor</b>	

Contact Hour	List of Topic Title (from content structure in Part A)	Topic # (from content structure in Part A)	Text/Ref Book/external resource
1-2	<ul style="list-style-type: none"> <li>• Introduction <ul style="list-style-type: none"> <li>o Information Vs Data Retrieval</li> <li>o Basic Concepts</li> <li>o The retrieval process</li> <li>o Taxonomy of IR</li> <li>o Classic IR and Alternative models</li> </ul> </li> </ul>	1.1, 1.2	R1 Ch1, Ch2
3-5	<ul style="list-style-type: none"> <li>• Boolean Retrieval <ul style="list-style-type: none"> <li>o Inverted index</li> <li>o Processing Boolean queries</li> <li>o Boolean Vs Ranked retrieval</li> <li>o Term vocabulary and postings lists</li> <li>o Phrase queries</li> </ul> </li> </ul>	2.1	T Ch2
6-7	<ul style="list-style-type: none"> <li>• Dictionary and Tolerant Retrieval <ul style="list-style-type: none"> <li>o Search Structures for dictionaries</li> <li>o Wildcard queries</li> <li>o Phonetic Correction</li> </ul> </li> </ul>	2.2	T Ch3
8-10	<ul style="list-style-type: none"> <li>• Index Construction and Compression <ul style="list-style-type: none"> <li>o Blocked sort-based Indexing</li> <li>o Single pass in-memory indexing</li> </ul> </li> </ul>	2.3	T Ch4

	<ul style="list-style-type: none"> <li>o Distributed and dynamic indexing</li> <li>o Dictionary comparison</li> <li>o Postings file compression</li> </ul>		
11-12	<ul style="list-style-type: none"> <li>• Vector Space Model <ul style="list-style-type: none"> <li>o Term frequency and weighting</li> <li>o The vector space model for scoring</li> <li>o Tf-idf functions</li> </ul> </li> </ul>	3.1, 3.2	T Ch6
13-15	<ul style="list-style-type: none"> <li>• Text Mining <ul style="list-style-type: none"> <li>o Classification <ul style="list-style-type: none"> <li>• Naïve Bayes</li> <li>• Vector space classification</li> <li>• Evaluating Classification</li> </ul> </li> </ul> </li> </ul>	4.1	T Ch13, 14
16-18	<ul style="list-style-type: none"> <li>• Text Mining <ul style="list-style-type: none"> <li>o Clustering <ul style="list-style-type: none"> <li>• Flat clustering</li> <li>• Hierarchical clustering</li> </ul> </li> </ul> </li> </ul>	4.2	T Ch16, 17
19	<ul style="list-style-type: none"> <li>• Web Search <ul style="list-style-type: none"> <li>o Web characteristics</li> <li>o The search user experience</li> <li>o Index size and estimation</li> </ul> </li> </ul>	5.1	T Ch19
20-22	<ul style="list-style-type: none"> <li>• Web Crawling and Indexes <ul style="list-style-type: none"> <li>o Crawling</li> <li>o Crawler Architecture</li> <li>o Distributed Indexes</li> </ul> </li> </ul>	5.2	T Ch20
23-24	<ul style="list-style-type: none"> <li>• Link Analysis <ul style="list-style-type: none"> <li>o The web as a graph</li> <li>o Google's page rank</li> <li>o Hub and Authorities (HITS)</li> </ul> </li> </ul>	5.3	T Ch21
25-27	<ul style="list-style-type: none"> <li>• Cross Lingual IR (CLIR) <ul style="list-style-type: none"> <li>o Language problems in IR</li> <li>o Translation Approaches</li> <li>o Handling Many Languages</li> <li>o Resources for CLIR</li> </ul> </li> </ul>	6.1, 6.2	R3 Ch2
28-29	<ul style="list-style-type: none"> <li>• Multimedia IR <ul style="list-style-type: none"> <li>o Basic Multimedia search technologies</li> <li>o Content Based Retrieval</li> </ul> </li> </ul>	7.1,7.2	R4 Ch2,3
30-31	<ul style="list-style-type: none"> <li>• Recommender System <ul style="list-style-type: none"> <li>o Collaborative recommendation</li> <li>o Content based recommendation</li> <li>o Other type &amp; hybrid recommendations</li> </ul> </li> </ul>	8.1	R2 Ch1-5
32	<ul style="list-style-type: none"> <li>• Review</li> </ul>		

**Work integration: Detailed plan**

No	Activity description (Examples are given below)
1	Apply Domain modelling concept to the work you are doing in the work place
2	Present the architecture of the software you are working on
3	Analyse the test plan of the software project you are working on and identify areas where it can be further improved
4	Seminar / talk by Project manager in the company on a topic of relevance to the course

**Evaluation Scheme**

Evaluation Component	Name (Quiz, Lab, Project, Mid term exam, End semester exam, etc)	Type (Open book, Closed book, Online, etc.)	Weight	Duration	Day, Date, Session, Time
EC - 1	Quiz	Online	10%		To be announced
	Assignment	Take home	10%		To be announced
EC - 2	Mid-Semester Test	Closed Book	30%	2 hrs	To be announced
EC - 3	Comprehensive Exam	Open Book	50%	3 hrs	To be announced

**Note** - Evaluation components can be tailored depending on the proposed model.

Syllabus for Mid-Semester Test (Closed Book): Topics in Weeks 1-7

Syllabus for Comprehensive Exam (Open Book): All topics given in plan of study

**Evaluation Guidelines:**

1. EC-1 consists of either two Assignments or three Quizzes. Announcements regarding the same will be made in a timely manner.
2. For Closed Book tests: No books or reference material of any kind will be permitted. Laptops/Mobiles of any kind are not allowed. Exchange of any material is not allowed.
3. For Open Book exams: Use of prescribed and reference text books, in original (not photocopies) is permitted. Class notes/slides as reference material in filed or bound form is permitted. However, loose sheets of paper will not be allowed. Use of calculators is permitted in all exams. Laptops/Mobiles of any kind are not allowed. Exchange of any material is not allowed.
4. If a student is unable to appear for the Regular Test/Exam due to genuine exigencies, the student should follow the procedure to apply for the Make-Up Test/Exam. The genuineness of the reason for absence in the Regular Exam shall be assessed prior to giving permission to appear for the Make-up Exam. Make-Up Test/Exam will be conducted only at selected exam centres on the dates to be announced later.

It shall be the responsibility of the individual student to be regular in maintaining the self-study schedule as given in the course handout, attend the lectures, and take all the prescribed evaluation components such as Assignment/Quiz, Mid-Semester Test and Comprehensive Exam according to the evaluation scheme provided in the handout.