

<b>Short Title:</b>	Text Analysis and Web Content Mining <b>APPROVED</b>
<b>Full Title:</b>	Text Analysis and Web Content Mining
<b>Module Code:</b>	ADSA H6014
<b>ECTS credits:</b>	10
<b>NFQ Level:</b>	9
<b>Module Delivered in</b>	<a href="#">1 programme(s)</a>
<b>Module Contributor:</b>	Markus Hofmann
<b>Module Description:</b>	Module aims include: • Investigate state of the art and research trends in text analysis and web content mining. • Critique and evaluate the performance of algorithms for both text analysis and web content mining
<b>Learning Outcomes:</b>	
<i>On successful completion of this module the learner will be able to</i>	
<ol style="list-style-type: none"> <li>1. Demonstrate an awareness and critical understanding of ways to extract key concepts and relationships from semi-structured and unstructured text, and structure them for data mining</li> <li>2. Discuss current research activities relating to text mining and web content mining</li> <li>3. Understand limitations of current information extraction techniques and the vision for the future</li> <li>4. Extract key concepts and relationships from semi-structured and unstructured data</li> <li>5. Apply prediction and clustering techniques to the prepared data, and critically evaluate the results in particular with the aid of modern text visualisation techniques.</li> <li>6. Independently research current trends and developments relating to the processing of semi-structured unstructured data</li> </ol>	

**Module Content & Assessment**

<b>Indicative Content</b>
<p><b>Preparing text documents for mining</b>          § Extracting key concepts, sentiments, and relationships from semi-structured and unstructured data; § Structural representations for text documents (e.g. Vector Space Model). § Apply appropriate visualisation techniques pre-model, model and post-model.</p>
<p><b>Mining the data</b>          § Learning methods for sparse, high dimensional data (e.g. support vector machines) and performance evaluation. § Clustering methods (e.g. k-means clustering; hierarchical clustering) and similarity measures for asymmetric data. § Visualisation techniques. § Case studies.</p>
<p><b>Web Crawling</b>          XPATH, web crawlers, regular expressions, crawling rules</p>
<p><b>Knowledge Extraction</b>          Concept extraction based on both syntactic and semantic natural language processing.</p>

<b>Indicative Assessment Breakdown</b>	<b>%</b>
Course Work Assessment %	100.00%

<b>Course Work Assessment %</b>				
<i>Assessment Type</i>	<i>Assessment Description</i>	<i>Outcome addressed</i>	<i>% of total</i>	<i>Assessment Date</i>
Reflective Journal	Students must prepare a portfolio of literary reviews and analysis covering a range of topics across all areas of the syllabus, and give an oral presentation of at least one of their research areas. For example, • An exploration of current trends in methods for structuring text in preparation for mining. • Investigation into at least one algorithm suitable for classifying text documents, and an analysis of current research in learning methods. • Investigation into at least one algorithm suitable for clustering text documents, and an analysis	1,2,3,6	30.00	Week 6
Practical/Skills Evaluation	Work through all stages of a text mining project life cycle using an appropriate text mining tool. For examples students would be presented with raw text and business objectives from which they would mine the data in an ongoing project during the semester with the following deliverables: • Step 1. Business Understanding: Evaluate the appropriate text mining function to be used to achieve the business objectives. • Step 2. Data Preparation: Structure the data in a format suitable for the relevant text mining function. • Step 3. Data Modelling: Apply the appropriate mining algorithm(s). • Evaluation: Evaluate the results.	4	40.00	Sem 1 End
Practical/Skills Evaluation	Students are asked to compile a unique data set by applying web crawling strategies and implementations. Once the data have been obtained, an appropriate analysis technique such as visualisation, classification, association rules or clustering need to be applied.	4	30.00	Sem 1 End

No Final Exam Assessment %
----------------------------

<b>Indicative Reassessment Requirement</b>
<p><b>Coursework Only</b>  <i>This module is reassessed solely on the basis of re-submitted coursework. There is no repeat written examination.</i></p>
<p><b>Reassessment Description</b>          As per course work</p>

ITB reserves the right to alter the nature and timings of assessment

**Indicative Module Workload & Resources**

<b>Resources</b>
<i>Recommended Book Resources</i>
<p>Hofmann, Chisholm 2016, <i>Text Mining and Visualization: Case Studies Using Open-Source Tools</i>, 1 Ed., Chapman &amp; Hall/CRC Data Mining and Knowledge Discovery Series [ISBN: 1482237571]</p> <p>Ashok Srivastava (Editor), Mehran Sahami (Editor), <i>Text Mining: Classification, Clustering, and Applications</i> [ISBN: 1420059408]</p> <p>Sholom M. Weiss... [et al.] 2005, <i>Text mining</i>, Springer New York [ISBN: 0-387-95433-3]</p> <p>editor, Michael W. Berry 2003, <i>Survey of text mining</i>, Springer New York [ISBN: 0-387-95563-1]</p>
<i>Supplementary Book Resources</i>
<p>Ian H. Witten, Alistair Moffat, Timothy C. Bell 1999, <i>Managing gigabytes</i>, Morgan Kaufmann Publishers San Francisco, Calif. [ISBN: 1558605703]</p>
<i>Recommended Article/Paper Resources</i>
<p>Madjid Tavana 2015, <i>International Journal of Knowledge Engineering and Data Mining</i> [ISSN: 1755-2087]  <a href="http://www.inderscience.com/jhome.php?jcode=ijkedm">http://www.inderscience.com/jhome.php?jcode=ijkedm</a></p>
<i>This module does not have any other resources</i>

**Module Delivered in**

Programme Code	Programme	Semester	Delivery
BN_KADSA_R	<a href="#">Master of Science in Computing in Applied Data Science &amp; Analytics</a>	3	Elective