

## Subject Description Form

<b>Subject Code</b>	COMP4433
<b>Subject Title</b>	Data Mining and Data Warehousing
<b>Credit Value</b>	3
<b>Level</b>	3
<b>Pre-requisite / Co-requisite / Exclusion</b>	<b>Pre-requisite:</b> COMP2411 or equivalent introductory database subject
<b>Objectives</b>	<p>This subject aims at equipping students with the latest knowledge and skills to:</p> <ol style="list-style-type: none"> <li>1. create a clean, consistent repository of data within a data warehouse for large corporations;</li> <li>2. utilise various techniques developed for data mining to discover interesting patterns in large databases;</li> <li>3. use existing commercial or public-domain tools to perform data mining tasks to solve real problems in business and commerce; and</li> <li>4. expose students to new techniques and ideas that can be used to improve the effectiveness of current data mining tools.</li> </ol>
<b>Intended Learning Outcomes</b>	<p>Upon completion of the subject, students will be able to:</p> <p><i>Professional/academic knowledge and skills</i></p> <ol style="list-style-type: none"> <li>(a) identify and analyse why there is a need for data warehouse in addition to traditional operational database systems, motivated by real examples;</li> <li>(b) conduct in-depth analysis of the key components in typical and advanced data warehouse architectures;</li> <li>(c) design a data warehouse and understand the process required to construct one;</li> <li>(d) identify and analyse why there is a need for data mining and in what ways it is different from traditional statistical techniques, motivated by real examples;</li> <li>(e) learn and master the algorithms made available by popular commercial data mining software;</li> <li>(f) solve real data mining problems by using the right tools to find interesting patterns;</li> <li>(g) obtain deep understanding of a typical knowledge discovery process;</li> <li>(h) obtain hands-on experience with some popular data mining software;</li> </ol> <p><i>Attributes for all-roundedness</i></p> <ol style="list-style-type: none"> <li>(i) apply data mining and data warehousing tools;</li> </ol>

	<p>(j) learn independently and search for relevant information to write reports to recommend appropriate data warehousing and data mining tools; and</p> <p>(k) generate innovative solutions individually or in groups and develop group work skills directly and indirectly.</p>										
<p><b>Subject Synopsis/ Indicative Syllabus</b></p>	<table border="1"> <thead> <tr> <th data-bbox="376 338 1473 412"><b>Topic</b></th> </tr> </thead> <tbody> <tr> <td data-bbox="376 412 1473 607"> <p><b>1. Introduction to Data Warehousing and Data Mining</b></p> <p>Introduction to data warehousing and data mining; possible application areas in business and finance; definitions and terminologies; types of data mining problems.</p> </td> </tr> <tr> <td data-bbox="376 607 1473 763"> <p><b>2. Data Warehousing</b></p> <p>Data warehouse and data warehousing; data warehouse and the industry; definitions; operational databases vs. data warehouses.</p> </td> </tr> <tr> <td data-bbox="376 763 1473 958"> <p><b>3. Data Warehouse Architecture and Design</b></p> <p>Data warehouse architecture and design; two-tier and three-tier architecture; star schema and snowflake schema; data characteristics; static and dynamic data; meta-data; data marts.</p> </td> </tr> <tr> <td data-bbox="376 958 1473 1153"> <p><b>4. Data Replication and Online Analytical Processing</b></p> <p>Data replication, data capturing and indexing, data transformation and cleansing; replicated data and derived data; Online Analytical Processing (OLAP); multidimensional databases; data cube.</p> </td> </tr> <tr> <td data-bbox="376 1153 1473 1310"> <p><b>5. Data Mining and Knowledge Discovery</b></p> <p>Data mining and knowledge discovery, the data mining lifecycle; pre-processing; data transformation; types of problems and applications.</p> </td> </tr> <tr> <td data-bbox="376 1310 1473 1467"> <p><b>6. Association Rules</b></p> <p>Mining of association rules; the Apriori algorithm; binary, quantitative and generalised association rules; interestingness measures.</p> </td> </tr> <tr> <td data-bbox="376 1467 1473 1662"> <p><b>7. Classification</b></p> <p>Classification; decision tree based algorithms; Bayesian approach; statistical approaches, nearest neighbour approach; neural network based approach; genetic algorithms based technique; evaluation of classification model.</p> </td> </tr> <tr> <td data-bbox="376 1662 1473 1818"> <p><b>8. Clustering</b></p> <p>Clustering; k-means algorithm; hierarchical algorithm; condorset; neural network and genetic algorithms based approach; evaluation of effectiveness.</p> </td> </tr> <tr> <td data-bbox="376 1818 1473 2031"> <p><b>9. Sequential Data Mining</b></p> <p>Sequential data mining; time dependent data and temporal data; time series analysis; sub-sequence matching; classification and clustering of temporal data; prediction.</p> </td> </tr> </tbody> </table>	<b>Topic</b>	<p><b>1. Introduction to Data Warehousing and Data Mining</b></p> <p>Introduction to data warehousing and data mining; possible application areas in business and finance; definitions and terminologies; types of data mining problems.</p>	<p><b>2. Data Warehousing</b></p> <p>Data warehouse and data warehousing; data warehouse and the industry; definitions; operational databases vs. data warehouses.</p>	<p><b>3. Data Warehouse Architecture and Design</b></p> <p>Data warehouse architecture and design; two-tier and three-tier architecture; star schema and snowflake schema; data characteristics; static and dynamic data; meta-data; data marts.</p>	<p><b>4. Data Replication and Online Analytical Processing</b></p> <p>Data replication, data capturing and indexing, data transformation and cleansing; replicated data and derived data; Online Analytical Processing (OLAP); multidimensional databases; data cube.</p>	<p><b>5. Data Mining and Knowledge Discovery</b></p> <p>Data mining and knowledge discovery, the data mining lifecycle; pre-processing; data transformation; types of problems and applications.</p>	<p><b>6. Association Rules</b></p> <p>Mining of association rules; the Apriori algorithm; binary, quantitative and generalised association rules; interestingness measures.</p>	<p><b>7. Classification</b></p> <p>Classification; decision tree based algorithms; Bayesian approach; statistical approaches, nearest neighbour approach; neural network based approach; genetic algorithms based technique; evaluation of classification model.</p>	<p><b>8. Clustering</b></p> <p>Clustering; k-means algorithm; hierarchical algorithm; condorset; neural network and genetic algorithms based approach; evaluation of effectiveness.</p>	<p><b>9. Sequential Data Mining</b></p> <p>Sequential data mining; time dependent data and temporal data; time series analysis; sub-sequence matching; classification and clustering of temporal data; prediction.</p>
<b>Topic</b>											
<p><b>1. Introduction to Data Warehousing and Data Mining</b></p> <p>Introduction to data warehousing and data mining; possible application areas in business and finance; definitions and terminologies; types of data mining problems.</p>											
<p><b>2. Data Warehousing</b></p> <p>Data warehouse and data warehousing; data warehouse and the industry; definitions; operational databases vs. data warehouses.</p>											
<p><b>3. Data Warehouse Architecture and Design</b></p> <p>Data warehouse architecture and design; two-tier and three-tier architecture; star schema and snowflake schema; data characteristics; static and dynamic data; meta-data; data marts.</p>											
<p><b>4. Data Replication and Online Analytical Processing</b></p> <p>Data replication, data capturing and indexing, data transformation and cleansing; replicated data and derived data; Online Analytical Processing (OLAP); multidimensional databases; data cube.</p>											
<p><b>5. Data Mining and Knowledge Discovery</b></p> <p>Data mining and knowledge discovery, the data mining lifecycle; pre-processing; data transformation; types of problems and applications.</p>											
<p><b>6. Association Rules</b></p> <p>Mining of association rules; the Apriori algorithm; binary, quantitative and generalised association rules; interestingness measures.</p>											
<p><b>7. Classification</b></p> <p>Classification; decision tree based algorithms; Bayesian approach; statistical approaches, nearest neighbour approach; neural network based approach; genetic algorithms based technique; evaluation of classification model.</p>											
<p><b>8. Clustering</b></p> <p>Clustering; k-means algorithm; hierarchical algorithm; condorset; neural network and genetic algorithms based approach; evaluation of effectiveness.</p>											
<p><b>9. Sequential Data Mining</b></p> <p>Sequential data mining; time dependent data and temporal data; time series analysis; sub-sequence matching; classification and clustering of temporal data; prediction.</p>											

	<p><b>10. Other Techniques</b></p> <p>Computation intelligence techniques; fuzzy logic, genetic algorithms and neural networks for data mining.</p> <p>Laboratory Experiment:</p> <table border="1" data-bbox="384 383 1465 607"> <tr> <td><b>Topic</b></td> </tr> <tr> <td>1. Discover Association rules and sequential patterns using data mining tools</td> </tr> <tr> <td>2. Discover Classification rules using data mining tools</td> </tr> <tr> <td>3. Discover Clusters using data mining tools</td> </tr> </table> <p>Case Study:</p> <ol style="list-style-type: none"> <li>1. Application of data mining techniques to solve real business problems.</li> <li>2. Attributes leading to success and failure of data warehousing projects tutorials when appropriate.</li> </ol>	<b>Topic</b>	1. Discover Association rules and sequential patterns using data mining tools	2. Discover Classification rules using data mining tools	3. Discover Clusters using data mining tools
<b>Topic</b>					
1. Discover Association rules and sequential patterns using data mining tools					
2. Discover Classification rules using data mining tools					
3. Discover Clusters using data mining tools					
<p><b>Teaching/ Learning Methodology</b></p>	<p>This subject consists mainly of class lectures and laboratory sessions. For the class lectures, various cases will be presented to help student understand why there is a need for data warehouse to be built and why data mining is important for modern day business intelligence. Students will be given time to participate in discussions when the cases are presented.</p> <p>All assignments and projects will also be given in the form of different cases collected so as to allow students to learn more about how data warehouse and data mining can be and have been used in real business environment. For the projects and assignments, students are expected to learn independently and think critically with minimise guidance. They are expected to practice their writing skills through project documentations and report writing. As students will work in teams on the project, they are expected to also learn to work with each other collaboratively.</p> <p>During laboratory sessions, students will be introduced to popular software products that can support the building of data warehouses and the mining of them. Students are expected to solve real data mining problems by using the right tools to find interesting patterns.</p>				

Assessment Methods in Alignment with Intended Learning Outcomes	Specific assessment methods/tasks	% weighting	Intended subject learning outcomes to be assessed										
			a	b	c	d	e	f	g	h	i	j	k
	<b>Continuous Assessment</b>	<b>55%</b>											
	1. Assignment		✓		✓	✓					✓	✓	
	2. Project					✓	✓	✓	✓	✓	✓	✓	✓
	<b>Examination</b>	<b>45%</b>	✓	✓	✓	✓	✓	✓	✓	✓	✓		
	<b>Total</b>	<b>100 %</b>											
			<p>The assessment consists of written assignments, a group project and an examination. For the assignments and projects, they are designed to ensure that students are able to achieve the learning outcomes intended for this subject. They are expected to tackle a number of cases drawn from different application areas in business and commerce so that they can understand why there is a need for data warehouse in addition to traditional operational database systems and why data mining is important for modern-day business intelligence. In addition, students will learn through the questions and cases, when a particular data warehouse architecture or when a particular data mining algorithm is useful and should be used. Questions in the assignments are expected to help students learning the details of the data mining algorithm and the use of popular data mining software. They are also expected to use such popular tool as Oracle Warehouse Builder to construct data warehouses. For the projects, students are expected to work in groups of three to four to tackle a real case involving the design of a data warehouse or the use of data mining to mine very large data bases. They are expected to learn how real-world problems in business and commerce should be tackled using real-world tools as Oracle's Warehouse Builder or IBM's Clementine data mining system. They are expected to learn independently and search for relevant information to write reports to recommend appropriate data warehousing and data mining tools. Students are expected to practice their writing skills with project document and report writing. They will learn to develop critical thinking and team work skills.</p>										
<b>Student Study Effort Expected</b>	Class contact:												
	▪ Lectures/Laboratory		39 Hrs.										
	▪ Tutorials		0 Hrs.										
	Other student study effort:												
	▪ Assignments and Case Studies		45 Hrs.										
	▪ Projects and Research		25 Hrs.										
	Total student study effort		109 Hrs.										

**Reading List  
and References**

**Reference Books:**

1. Han, Jiawei and Kamber, Micheline, *Data Mining: Concepts and Techniques*, 3<sup>rd</sup> Edition, Morgan Kaufmann, 2012.
2. Golfarelli, Matteo and Rizzi, Stefano, *Data Warehouse Design: Modern Principles and Methodologies*, McGraw-Hill, 2009.
3. Inmon, W.H., Strauss, Derek and Neushloss, Genia, *DW 2.0: The Architecture for the Next Generation of Data Warehousing*, Morgan Kaufmann, 2008.
4. Rokach, Lior and Maimon, Oded Z., *Data Mining with Decision Trees: Theory and Applications*, World Scientific, 2008.
5. Witten, Ian H., Frank, Eibe and Hall, Mark A., *Data Mining: Practical Machine Learning Tools and Techniques*, 3<sup>rd</sup> Edition, Morgan Kaufmann, 2011.
6. Westphal, Christopher, *Data Mining for Intelligence, Fraud & Criminal Detection: Advanced Analytics & Information Sharing Technologies*, CRC Press, 2008.
7. Cox, Earl, *Fuzzy Modeling and Genetic Algorithms for Data Mining and Exploration*, Morgan Kaufmann, 2005.
8. Liu, Bing, *Web Data Mining: Exploring Hyperlinks, Contents, and Usage Data*, Springer, Berlin Heidelberg, 2009.
9. Tsitsis, Konstantinos K. and Chorianopoulos, Antonios, *Data Mining Techniques in CRM: Inside Customer Segmentation*, Wiley, 2010.
10. Shapiro, A.F. and Jain, L.C., *Intelligent and Other Computational Techniques in Insurance: Theory and Applications*, World Scientific, 2003.