

Subject Description Form

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

	<p>data warehousing tools;</p> <p>(j) learn independently and search for relevant information to write reports to recommend appropriate data warehousing and data mining tools.</p> <p>(k) Solve complex problems individually or in groups and develop group work skills directly and indirectly.</p>													
<p>Subject Synopsis/ Indicative Syllabus</p>	<table border="1"> <thead> <tr> <th data-bbox="432 456 1445 495">Topic</th> </tr> </thead> <tbody> <tr> <td data-bbox="432 495 1445 645"> <p>1. Introduction to data warehousing and data mining 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="432 645 1445 757"> <p>2. Data warehousing Data warehouse and data warehousing; data warehouse and the industry; definitions; operational databases vs. data warehouses.</p> </td> </tr> <tr> <td data-bbox="432 757 1445 907"> <p>3. Data warehouse architecture and design 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="432 907 1445 1057"> <p>4. Data Replication and Online Analytical Processing 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="432 1057 1445 1169"> <p>5. Data mining and knowledge discovery 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="432 1169 1445 1281"> <p>6. Association rules Mining of association rules; the Apriori algorithm; binary, quantitative and generalized association rules; interestingness measures.</p> </td> </tr> <tr> <td data-bbox="432 1281 1445 1476"> <p>7. Classification Classification; decision tree based algorithms; Bayesian approach; statistical approaches, nearest neighbor approach; neural network based approach; genetic algorithms based technique; evaluation of classification model.</p> </td> </tr> <tr> <td data-bbox="432 1476 1445 1626"> <p>8. Clustering Clustering; k-means algorithm; hierarchical algorithm; Condorset; neural network and genetic algorithms based approach; evaluation of effectiveness.</p> </td> </tr> <tr> <td data-bbox="432 1626 1445 1776"> <p>9. Sequential data mining 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> <tr> <td data-bbox="432 1776 1445 1888"> <p>10. Other techniques Computation intelligence techniques; fuzzy logic, genetic algorithms and neural networks for data mining.</p> </td> </tr> </tbody> </table> <p>Laboratory Experiment:</p> <table border="1"> <thead> <tr> <th data-bbox="432 2029 1270 2089">Topic</th> </tr> </thead> <tbody> <tr> <td data-bbox="432 2089 1270 2089"></td> </tr> </tbody> </table>	Topic	<p>1. Introduction to data warehousing and data mining Introduction to data warehousing and data mining; possible application areas in business and finance; definitions and terminologies; types of data mining problems.</p>	<p>2. Data warehousing Data warehouse and data warehousing; data warehouse and the industry; definitions; operational databases vs. data warehouses.</p>	<p>3. Data warehouse architecture and design 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>4. Data Replication and Online Analytical Processing 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>5. Data mining and knowledge discovery Data mining and knowledge discovery, the data mining lifecycle; pre-processing; data transformation; types of problems and applications.</p>	<p>6. Association rules Mining of association rules; the Apriori algorithm; binary, quantitative and generalized association rules; interestingness measures.</p>	<p>7. Classification Classification; decision tree based algorithms; Bayesian approach; statistical approaches, nearest neighbor approach; neural network based approach; genetic algorithms based technique; evaluation of classification model.</p>	<p>8. Clustering Clustering; k-means algorithm; hierarchical algorithm; Condorset; neural network and genetic algorithms based approach; evaluation of effectiveness.</p>	<p>9. Sequential data mining Sequential data mining; time dependent data and temporal data; time series analysis; sub-sequence matching; classification and clustering of temporal data; prediction.</p>	<p>10. Other techniques Computation intelligence techniques; fuzzy logic, genetic algorithms and neural networks for data mining.</p>	Topic	
Topic														
<p>1. Introduction to data warehousing and data mining Introduction to data warehousing and data mining; possible application areas in business and finance; definitions and terminologies; types of data mining problems.</p>														
<p>2. Data warehousing Data warehouse and data warehousing; data warehouse and the industry; definitions; operational databases vs. data warehouses.</p>														
<p>3. Data warehouse architecture and design 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>4. Data Replication and Online Analytical Processing 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>5. Data mining and knowledge discovery Data mining and knowledge discovery, the data mining lifecycle; pre-processing; data transformation; types of problems and applications.</p>														
<p>6. Association rules Mining of association rules; the Apriori algorithm; binary, quantitative and generalized association rules; interestingness measures.</p>														
<p>7. Classification Classification; decision tree based algorithms; Bayesian approach; statistical approaches, nearest neighbor approach; neural network based approach; genetic algorithms based technique; evaluation of classification model.</p>														
<p>8. Clustering Clustering; k-means algorithm; hierarchical algorithm; Condorset; neural network and genetic algorithms based approach; evaluation of effectiveness.</p>														
<p>9. Sequential data mining Sequential data mining; time dependent data and temporal data; time series analysis; sub-sequence matching; classification and clustering of temporal data; prediction.</p>														
<p>10. Other techniques Computation intelligence techniques; fuzzy logic, genetic algorithms and neural networks for data mining.</p>														
Topic														

1. Knowledge discovery lifecycle using CRISP-DM
2. Discover Association rules and sequential patterns using Clementine
3. Discover Classification rules using Clementine
4. Discover Clusters using Clementine

Case Study:

- Application of data mining techniques to solve real business problems.
- Attributes leading to success and failure of data warehousing projects tutorials when appropriate.

Teaching/Learning Methodology

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.

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 minimize 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.

During laboratory sessions, students will be introduced to popular software products from Oracle and IBM 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.

Assessment Methods in Alignment with Intended Learning Outcomes

Specific assessment methods/tasks	% weighting	Intended subject learning outcomes to be assessed (Please tick as appropriate)										
		a	b	C	d	e	f	g	h	i	j	k
1. Assignments	55%	✓		✓	✓					✓	✓	
2. Project					✓	✓	✓	✓	✓	✓	✓	✓
3. Examination	45%	✓	✓	✓	✓	✓	✓	✓	✓	✓		
Total	100 %											

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

	<p>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>	
<p>Student Study Effort Expected</p>	<p>Class contact:</p>	
	<ul style="list-style-type: none"> ▪ Lecture/Laboratory 	<p>39 Hrs.</p>
	<ul style="list-style-type: none"> ▪ Tutorial 	<p>0 Hrs.</p>
	<p>Other student study effort:</p>	
	<ul style="list-style-type: none"> ▪ Assignments and case studies 	<p>45 Hrs.</p>
	<ul style="list-style-type: none"> ▪ Projects and research 	<p>25 Hrs.</p>
	<p>Total student study effort</p>	<p>109 Hrs.</p>
<p>Reading List and References</p>	<p>Reference Books:</p> <ol style="list-style-type: none"> 1. Chan, K.C.C., <i>Course Notes and Lab Manuals for COMP417</i>, 2009. 2. Inmon, W.H., Strauss, D., and Neushloss, G., <i>DW 2.0: The Architecture for the Next Generation of Data Warehousing</i>, Morgan Kaufmann, 2008. 3. Golfarelli, M., and Rizzi, S., <i>Data Warehouse Design: Modern Principles and Methodologies</i>, McGraw-Hill, 2009. 4. Rokach, L., and Maimon, O., <i>Data Mining with Decision Trees: Theory and Applications</i>, World Scientific, 2008. 5. Witten, I.H., Frank, E., Hall, M.A., <i>Data Mining, Third Edition: Practical Machine Learning Tools and Techniques</i>, Morgan Kaufmann, 2011. 6. Westphal., C., <i>Data Mining for Intelligence, Fraud & Criminal Detection: Advanced Analytics & Information Sharing Technologies</i>, CRC Press, 2008. 7. Cox, E., <i>Fuzzy Modeling and Genetic Algorithms for Data Mining and Exploration</i>, Morgan Kaufmann, 2005. 8. Liu, B., <i>Web Data Mining: Exploring Hyperlinks, Contents, and Usage Data</i>, Springer, Berlin Heidelberg, 2009. 9. Tsiptsis, K., and Chorianopoulos, A., <i>Data Mining Techniques in CRM: Inside Customer Segmentation</i>, Wiley, 2010. 10. Han, J. and Kamber, M., <i>Data Mining: Concepts and Techniques</i>, 2nd Edition, Morgan Kaufmann, 2005. 11. Shapiro, A.F., and Jain, L.C., <i>Intelligent and Other Computational Techniques in Insurance: Theory and Applications</i>, World Scientific, 2003. 	