

## Subject Description Form

<b>Subject Code</b>	COMP4434
<b>Subject Title</b>	Big Data Analytics
<b>Credit Value</b>	3
<b>Level</b>	4
<b>Pre-requisites</b>	AMA1104 Introductory Probability COMP1011 Programming Fundamentals COMP2011 Data Structures COMP2411 Database Systems
<b>Objectives</b>	<p>The objectives of this subject are to:</p> <ol style="list-style-type: none"> <li>1. introduce students the concept and challenge of big data (3 V's: volume, velocity, and variety).</li> <li>2. teach students in applying skills and tools to manage and analyze the big data.</li> </ol>
<b>Intended Learning Outcomes</b>	<p>Upon completion of the subject, students will be able to:</p> <ol style="list-style-type: none"> <li>(a) understand the concept and challenge of big data and why existing technology is inadequate to analyze the big data;</li> <li>(b) collect, manage, store, query, and analyze various form of big data; and</li> <li>(c) gain hands-on experience on large-scale analytics tools to solve some open big data problems; and</li> <li>(d) understand the impact of big data for business decisions and strategy.</li> </ol>
<b>Subject Synopsis/ Indicative Syllabus</b>	<ol style="list-style-type: none"> <li>1. Introduction to Big Data: The 3 V's, their challenges and application domains.</li> <li>2. Collection of Big Data: Eventual Consistency and NoSQL systems MongoDB, Google BigTable</li> <li>3. Large-Scale Data Analytics Systems: Auto-Parallel Data Programming; MapReduce, Hive, and Parallel Databases</li> <li>4. Basic Statistical Analysis: Fruad and Benfords Law, Bayesian Introduction, Heteroskedasticity</li> <li>5. Machine Learning Systems for Big Data</li> <li>6. Graph Analytics: Graph structures (diameter, connectivity, centrality), PageRank, Triangle counting</li> <li>7. Sentiment Analysis</li> </ol>

	8. Data Visualization: Data types and dimensions; Visual encoding and perception						
<b>Teaching/Learning Methodology</b>	A mix of lectures and lab sessions is used to deliver the various topics in this subject. Lectures are conducted to initiate students with the concepts and techniques of big data. Students are given the opportunity to gain hands-on experience on both open-source and commercial big data analytics software during the laboratory sessions.						
<b>Assessment Methods in Alignment with Intended Learning Outcomes</b>	Specific assessment methods/tasks	% weighting	Intended subject learning outcomes to be assessed (Please tick as appropriate)				
			a	b	c	d	
	1. Lab exercises/Assignments	60%	x	x	x	x	
	2. Project		x	x	x	x	
	3. Quiz		x	x			
	4. Examination	40%	x	x		x	
	Total	100 %					
<p>Explanation of the appropriateness of the assessment methods in assessing the intended learning outcomes:</p> <p>Continuous assessments consist of a project, assignments, lab exercises, and quizzes, which are designed to facilitate students to achieve intended learning outcomes. Lab exercise is designed to encourage students to acquire deep understanding of the relevant knowledge, practice in order to enrich their hands-on experience with various software tools. The project is designed to enhance students' ability to acquire the understanding and using different knowledge, principles, techniques, tools to solve a real problem through team. Quizzes are to ensure the students understand the concepts.</p> <p>Examination will evaluate student's understanding and usage of big data technologies.</p>							
<b>Student Study Effort Expected</b>	Class contact:						
	▪ Lecture						26 Hrs.
	▪ Tutorial/Laboratory						13 Hrs.
	Other student study effort:						

	<ul style="list-style-type: none"> <li>▪ Review the lecture</li> </ul>	28 Hrs.
	<ul style="list-style-type: none"> <li>▪ Review the lab</li> </ul>	14 Hrs.
	<ul style="list-style-type: none"> <li>▪ Work on the project</li> </ul>	15 Hrs.
	<ul style="list-style-type: none"> <li>▪ Prepare the quizzes</li> </ul>	9 Hrs.
	<ul style="list-style-type: none"> <li>▪ Prepare the examination</li> </ul>	11 Hrs.
	Total student study effort	116 Hrs.
Reading List and References	<ol style="list-style-type: none"> <li>1. How Vertica Was the Star of the Obama Campaign, and Other Revelations</li> <li>2. Cohen et al. "MAD Skills: New Analysis Practices for Big Data", 2009</li> <li>3. Ullman, Rajaraman, Mining of Massive Datasets, Chapter 2</li> <li>4. Stonebraker et al., "MapReduce and Parallel DBMS's: Friends or Foes?", Communications of the ACM, January 2010.</li> <li>5. Dean and Ghemawat, "MapReduce: A Flexible Data Processing Tool", Communications of the ACM, January 2010.</li> <li>6. Rick Cattell, "Scalable SQL and NoSQL Data Stores", SIGMOD Record, December 2010 (39:4)</li> <li>7. Elmagarmid, et. al. "Duplicate Record Detection: A Survey"</li> <li>8. Koudas, et. al. "Record Linkage: Similarity Measures and Algorithms"</li> <li>9. Chapter 3 of A Handbook of Statistical Analyses Using R</li> <li>10. Gregory Park on overfitting to the leaderboard in a Kaggle Competition</li> <li>11. Xindong Wu et al., Top 10 Algorithms in Data Mining, Knowledge and Information Systems, 14(2008), 1: 1-37. (read C4.5)</li> <li>12. Ullman, Rajaraman, Mining of Massive Datasets , Chapter 1</li> <li>13. Pedro Domingos, A Few Useful Things to Know about Machine Learning, CACM 55(10), 2012</li> <li>14. Pat Hanaran, Tools for Data Enthusiasts</li> <li>15. Jeffrey Heer, Michael Bostock, Vadim Ogievetsky, A Tour through the Visualization Zoo, Communications of the ACM, Volume 53 Issue 6, June 2010</li> <li>16. Howard Wen, "Big Ethics for Big Data", O'Reilly Media</li> </ol>	