

## Subject Description Form

<b>Subject Code</b>	COMP5434
<b>Subject Title</b>	Big Data Computing
<b>Credit Value</b>	3
<b>Level</b>	5
<b>Pre-requisites</b>	Knowledge in database systems, machine learning and data analytics is preferred.
<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;</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 traditional technology is inadequate to analyze the big data;</li> <li>b. understand how to collect, manage, store, and query various form of big data;</li> <li>c. familiar with the classical data analysis and machine learning algorithms;</li> <li>d. familiar with large-scale analytics tools to solve some open big data problems; and</li> <li>e. analyze the impact of big data for real-world business decisions and strategy.</li> </ol>
<b>Subject Synopsis/ Indicative Syllabus</b>	<ol style="list-style-type: none"> <li>1. Introduction to Big Data: Different V's, their challenges and application domains.</li> <li>2. Cloud Computing Basics: Software as a service (SaaS), Platform as a Service (PaaS), Infrastructure as a Service (IaaS), Desktop as a Service (DaaS), Public, Private and Enterprise Cloud.</li> <li>3. Big Data Computing: Concepts, Platform, Service, and Tools</li> <li>4. Large-Scale Programming Abstraction: MapReduce and its open source implementation of Hadoop</li> <li>5. Large-Scale Data Processing Framework: Apache Spark and its Built-in Modules</li> <li>6. Large-Scale Database Management: NoSQL and other tools, e.g. MongoDB, Google BigTable, etc.</li> <li>7. Machine Learning Systems for Big Data: Methods and Tools</li> <li>8. Big Data Visualization: Data types and dimensions; Visual encoding and perception</li> <li>9. Big Data Case Studies</li> </ol>

<b>Teaching/Learning Methodology</b>	<p>A mix of lectures, discussions and case studies.</p> <p>Class activities include lectures, tutorials, laboratory works and seminars.</p>																																																		
<b>Assessment Methods in Alignment with Intended Learning Outcomes</b>	<table border="1" data-bbox="534 358 1479 900"> <thead> <tr> <th data-bbox="534 358 810 526" rowspan="2">Specific assessment methods/tasks</th> <th data-bbox="815 358 965 526" rowspan="2">% weighting</th> <th colspan="5" data-bbox="970 358 1479 459">Intended subject learning outcomes to be assessed (Please tick as appropriate)</th> </tr> <tr> <th data-bbox="970 465 1062 526">a</th> <th data-bbox="1067 465 1160 526">b</th> <th data-bbox="1165 465 1257 526">c</th> <th data-bbox="1262 465 1355 526">d</th> <th data-bbox="1359 465 1479 526">e</th> </tr> </thead> <tbody> <tr> <td data-bbox="534 533 810 622">1. Assignments or lab works</td> <td data-bbox="815 533 965 757" rowspan="3">55</td> <td data-bbox="970 533 1062 622">✓</td> <td data-bbox="1067 533 1160 622">✓</td> <td data-bbox="1165 533 1257 622">✓</td> <td data-bbox="1262 533 1355 622">✓</td> <td data-bbox="1359 533 1479 622">✓</td> </tr> <tr> <td data-bbox="534 629 810 696">2. Project</td> <td data-bbox="970 629 1062 696">✓</td> <td data-bbox="1067 629 1160 696">✓</td> <td data-bbox="1165 629 1257 696">✓</td> <td data-bbox="1262 629 1355 696">✓</td> <td data-bbox="1359 629 1479 696">✓</td> </tr> <tr> <td data-bbox="534 703 810 770">3. Quiz</td> <td data-bbox="970 703 1062 770">✓</td> <td data-bbox="1067 703 1160 770">✓</td> <td data-bbox="1165 703 1257 770">✓</td> <td data-bbox="1262 703 1355 770">✓</td> <td data-bbox="1359 703 1479 770"></td> </tr> <tr> <td data-bbox="534 777 810 844">4. Examination</td> <td data-bbox="815 777 965 844">45</td> <td data-bbox="970 777 1062 844">✓</td> <td data-bbox="1067 777 1160 844">✓</td> <td data-bbox="1165 777 1257 844">✓</td> <td data-bbox="1262 777 1355 844"></td> <td data-bbox="1359 777 1479 844">✓</td> </tr> <tr> <td data-bbox="534 851 810 900">Total</td> <td data-bbox="815 851 965 900">100</td> <td colspan="5" data-bbox="970 851 1479 900"></td> </tr> </tbody> </table> <p data-bbox="534 952 1460 1019">Explanation of the appropriateness of the assessment methods in assessing the intended learning outcomes:</p> <p data-bbox="534 1041 1460 1310">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 good 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 data-bbox="534 1332 1396 1400">Examination will evaluate student's understanding and usage of big data technologies.</p>						Specific assessment methods/tasks	% weighting	Intended subject learning outcomes to be assessed (Please tick as appropriate)					a	b	c	d	e	1. Assignments or lab works	55	✓	✓	✓	✓	✓	2. Project	✓	✓	✓	✓	✓	3. Quiz	✓	✓	✓	✓		4. Examination	45	✓	✓	✓		✓	Total	100					
Specific assessment methods/tasks	% weighting	Intended subject learning outcomes to be assessed (Please tick as appropriate)																																																	
		a	b	c	d	e																																													
1. Assignments or lab works	55	✓	✓	✓	✓	✓																																													
2. Project		✓	✓	✓	✓	✓																																													
3. Quiz		✓	✓	✓	✓																																														
4. Examination	45	✓	✓	✓		✓																																													
Total	100																																																		
<b>Student Study Effort Expected</b>	<p>Class contact:</p> <p>Class activities (lecture, tutorial, lab, etc.)</p> <p>Other student study effort:</p> <p>Assignments, Quizzes, Projects, Examination</p> <p>Total student study effort</p>					<p></p> <p>39 Hrs.</p> <p></p> <p>65 Hrs.</p> <p><b>104 Hrs.</b></p>																																													
<b>Reading List and References</b>	<ol style="list-style-type: none"> <li>1. Jared Dean, Big Data, Data Mining, and Machine Learning: Value Creation for Business Leaders and Practitioners. Wiley, 2014.</li> <li>2. Steele, Julie, and Noah Iliinsky, Beautiful visualization: looking at data through the eyes of experts, O'Reilly Media, Inc., 2010.</li> <li>3. Dean, Jeffrey and Ghemawat, Sanjay, "MapReduce: simplified data processing on large clusters", Communications of the ACM, January 2008.</li> <li>4. Stonebraker, M., Abadi, D., DeWitt, David J., Madden, S., Paulson, E., Pavlo, A. and Rasin, A., "MapReduce and Parallel DBMS's: Friends or Foes?", Communications of the ACM, January 2010.</li> </ol>																																																		

5. Dean, Jeffrey and Ghemawat, Sanjay, "MapReduce: A Flexible Data Processing Tool", Communications of the ACM, January 2010.
6. Lin, Jimmy and Dyer, Chris, Data-Intensive Text Processing with MapReduce, Morgan and Claypool, 2010.
7. K. Shvachko, H. Kuang, S. Radia and R. Chansler, "The Hadoop Distributed File System", IEEE Symposium on Mass Storage Systems and Technologies, 2010.
8. White, Tom, Hadoop: The definitive guide, O'Reilly Media, Inc., 2012.
9. Cattell, Rick, "Scalable SQL and NoSQL Data Stores", ACM SIGMOD Record, Volume 39, Issue 4, December 2010.
10. Chodorow, Kristina. MongoDB: the definitive guide: powerful and scalable data storage, O'Reilly Media, Inc., 2013.
11. Silberschatz, Abraham, Henry F. Korth, and Shashank Sudarshan, Database System Concepts, 7th Edition, 2019.
12. Page, Lawrence and Brin, Sergey and Motwani, Rajeev and Winograd, Terry, "The PageRank Citation Ranking: Bringing Order to the Web", Technical Report, Stanford InfoLab, 1999.
13. Wu, X.D., Kumar, V., Quinlan, J. Ross, Ghosh, J., Yang, Q. et al., "Top 10 Algorithms in Data Mining, Knowledge and Information Systems", Journal of Knowledge and Information Systems, Volume 14, Issue 1, page 1-37, 2007.
14. Leskovec, Rajaraman, Ullman, Mining of Massive Datasets, 2nd Edition, Cambridge University Press, 2014.
15. Tan, Pang-Ning, Michael Steinbach, and Vipin Kumar, Introduction to data mining, Pearson Education India, 2016.
16. Hastie, Trevor, Robert Tibshirani, and Jerome Friedman, The Elements of Statistical Learning: Data mining, Inference, and Prediction, Springer Science & Business Media, 2009.
17. Bishop, Christopher M., Pattern Recognition and Machine Learning, Springer, 2006.
18. Goodfellow, Ian, et al., Deep Learning: Adaptive Computation and Machine Learning series, MIT press, 2016.
19. McKinney, W., Python for data analysis: Data wrangling with Pandas, NumPy, and IPython, O'Reilly Media, Inc., 2012.
20. Hothorn, Torsten and Everitt, Brian S., A Handbook of Statistical Analyses Using R, CRC Press, 2014.
21. Géron, A., Hands-on machine learning with Scikit-Learn, Keras, and TensorFlow: Concepts, tools, and techniques to build intelligent systems, O'Reilly Media, 2019.
22. Nickoloff, J., Docker in action, Manning Publications Co., 2016.