

## Coursework Specification

### CW\_Specification\_CSI\_6\_DMA\_23/24

Read this coursework specification carefully, it tells you how you are going to be assessed, how to submit your coursework on-time and how (and when) you'll receive your marks and feedback.

<b>Module Code</b>	CSI_6_DMA
<b>Module Title</b>	Data Mining and Big Data Analytics
<b>Lecturer</b>	Daqing Chen; George Bamfo
<b>% of Module Mark</b>	60%
<b>Distributed</b>	04/10/2023
<b>Submission Method</b>	Submit online via this Module's Moodle site
<b>Submission Deadline</b>	17:00, 11/12/2023
<b>Release of Feedback &amp; Marks</b>	Feedback and provisional marks will be available in Grades on Moodle from 03/01/2024

### Coursework Aim:

This individual coursework project is about using analytics to address real-world problems. The aim of this coursework is to evaluate your understanding of the basic theories, concepts, methodologies, and the various algorithms in data mining, and your skills of using Python for analytics tools.

### Coursework Details:

<b>Type:</b>	Project report
<b>Overall View:</b>	<p>This assignment is to be undertaken <b>individually</b>. You should plan your work thoroughly and have regular discussions with your module tutor to resolve any issues you may have during this project.</p> <p>The assignment involves analysing a real-world dataset and identifying its underlying structural patterns and models using appropriate data mining techniques and algorithms. These patterns and models are intended to be used to address certain business concerns. A dataset will be assigned to you by your tutor.</p>

	<p>You are expected to produce a written report for this data mining project.</p>
<b>Tasks:</b>	<p>You are required to undertake the following tasks in this project:</p> <ol style="list-style-type: none"> <li>1. <b><u>Business Understanding</u></b> <ul style="list-style-type: none"> <li>• Download the dataset assigned to you from the module Moodle site along with the data description file.</li> <li>• Read the data description file to learn the nature of the data, such as what is the data about, where it comes from, which certain business context it is associated with, etc.</li> <li>• Examine the dataset within its business context to identify meaningful problems that potentially can be addressed by using analytics.</li> <li>• Translate the business problems to appropriate data mining problems and tasks.</li> <li>• You may also refer to any articles published relevant to the dataset.</li> </ul> </li> <li>2. <b><u>Data Understanding</u></b> <ul style="list-style-type: none"> <li>• Perform initial data exploration to get to know more about the dataset, such as the total number of instances in the dataset, the number of attributes (variables), the data type of each attribute, and the basic statistics of each attribute, including value range, average, standard deviation, skewness, kurtosis, and mode, etc.</li> <li>• Identify any data quality issues, including missing values, outliers, extreme values, incomparable value ranges of variables, and imbalanced classes, etc.</li> <li>• Determine if the dataset is appropriate to be used for addressing the business problems identified in Task 1. If not, re-do Task 1.</li> </ul> </li> <li>3. <b><u>Data Preparation</u></b> <ul style="list-style-type: none"> <li>• Choose appropriate methods for data pre-processing, which includes dealing with missing values, tackling outliers, extreme values, and imbalanced classes, changing data types, reducing dimensionality, and conducting data transformation and normalisation, etc., wherever appropriate.</li> <li>• Identify correlations among certain variables.</li> </ul> </li> </ol>

	<ul style="list-style-type: none"> <li>• Determine which and how each attribute should be used in your analysis.</li> <li>• Divide the whole dataset into several subsets to be used for training, test and validation in predictive modelling.</li> </ul> <p>4. <b><u>Modelling</u></b></p> <ul style="list-style-type: none"> <li>• Use the pre-processed dataset to perform the data mining tasks you have identified in Task 1.</li> <li>• Choose appropriate techniques and algorithms for your analysis: Choose either <i>k</i>-means clustering or association rule analysis for descriptive modelling, and choose either decision tree or regression for predictive modelling.</li> <li>• Determine appropriate settings of the algorithms to be applied, e.g., how many clusters to use in <i>k</i>-means clustering.</li> <li>• Re-do data preparation in Task 3 if needed.</li> </ul> <p>5. <b><u>Evaluation</u></b></p> <ul style="list-style-type: none"> <li>• Provide an explicit and concise description and explanation of the descriptive and predictive models you have created. Examine and explain what patterns and insight have been identified.</li> <li>• Evaluate the performance of the predictive models in terms of evaluate the performance of the predictive models in terms of various measures applicable, such as accuracy, SSE (sum of squared errors), generalisation ability, simplicity and cost etc.</li> <li>• Discuss how the descriptive and predictive models created can be used to address the original business problems identified in Task 1.</li> <li>• Summarise your main findings from the project.</li> </ul>
<b>Word Count:</b>	<p>As a guide, aim for 2500-3000 words. The maximum word limit is 3000 words. If the total word limit is exceeded, it will affect the marks awarded to the project presentation.</p> <p>Footnotes will not count towards word count totals but must only be used for referencing, not for the provision of additional text. The bibliography will not count towards the word total.</p>

<b>Presentation:</b>	<ul style="list-style-type: none"> <li>• Report must contain Title page, Table of Contents, Abstract, Conclusion, and References.</li> <li>• Work must be referenced, and a bibliography provided.</li> <li>• Work must be submitted as a Word document (.doc/docx) or a PDF.</li> <li>• Course work must be submitted using Arial font size 11 (or larger if you need to), with a minimum of 1.5 line spacing.</li> <li>• Your student number must appear at the front of the coursework. Your name must <b><u>not</u></b> be on your coursework.</li> </ul>
<b>Referencing:</b>	Harvard Referencing should be used, see your <a href="#">Library Subject Guide</a> for guides and tips on referencing.
<b>Regulations:</b>	<p>Make sure you understand the <a href="#">University Regulations</a> on expected academic practice and academic misconduct. Note in particular:</p> <ul style="list-style-type: none"> <li>▪ Your work must be your own. Markers will be attentive to both the plausibility of the sources provided as well as the consistency and approach to writing of the work. Simply, if you do the research and reading, and then write it up on your own, giving the reference to sources, you will approach the work in the appropriate way and will cause not give markers reason to question the authenticity of the work.</li> <li>▪ All quotations must be credited and properly referenced. Paraphrasing is still regarded as plagiarism if you fail to acknowledge the source for the ideas being expressed.</li> </ul> <p><b>TURNITIN:</b> When you upload your work to the Moodle site it will be checked by anti-plagiarism software.</p>

## Learning Outcomes

This coursework will partially assess the following learning outcomes for this module as indicated by \*.

### Knowledge and Understanding

On successful completion of this module, you will be able to:

- Describe and explain the concepts of data mining including the techniques and algorithms for problem solving and creating competitive advantage. \*

### **Intellectual Skills**

On successful completion of this module, you will be able to:

- Critically evaluate different types of data mining tasks in relation to various business and scientific problems, including descriptive modelling and predictive modelling, including cluster analysis, association analysis, and decision and regression for classification and prediction. \*

### **Practical Skills**

On successful completion of this module, you will be able to:

- Transfer a business and/or scientific problem into an appropriate data mining problem. \*
- Creatively apply data mining tools and platforms such as Python package.\*

### **Transferable Skills**

On successful completion of this module, you will be able to:

- Analyse and develop solutions for a wide range of business and scientific problems. \*

## **Assessment Criteria and Weighting**

LSBU marking criteria have been developed to help tutors give you clear and helpful feedback on your work. They will be applied to your work to help you understand what you have accomplished, how any mark given was arrived at, and how you can improve your work in future.

	Criteria	Feedforward comments						
		100 - 80%	79 - 70%	69 - 60%	59 - 50%	49 - 40%	39 - 30%	29 - 0%
10%	<b>1. Business Understanding</b>	Exceptionally thorough and clear analysis of business concerns and associated data mining tasks.	Thorough and clear analysis of business concerns and associated data mining tasks.	Clear analysis of business concerns and associated data mining tasks to a certain depth.	Clear analysis of business concerns and associated data mining tasks. Probably lack some in-depth view.	Basic analysis of the key business concerns and associated data mining tasks.	Inadequate analysis of business concerns and associated data mining tasks. Lack clarity and relevance.	Little or no analysis of business concerns and associated data mining tasks.
10%	<b>2. Data Understanding</b>	Exceptionally excellent and creative initial data exploration with effective means. Thorough summary of the dataset. Excellent analysis of data quality issues and the role of each attribute. Excellent use of Python.	Excellent initial data exploration with effective means. Thorough summary of the dataset. Excellent analysis of data quality issues and the role of each attribute. Excellent use of Python.	Good initial data exploration performed with appropriate means. Clear summary of the dataset. Good analysis of data quality issues and the role of each attribute. Good and flexible use of Python.	Essential initial data exploration performed. Essential analysis of data quality issues and the role of each attribute. Good use of Python.	Limited simple initial data exploration. Probably lack some relevance and/or clarity. Limited use of Python.	Inadequate and/or inappropriate initial data exploration performed. Lack clarity and relevance. Poor use of Python.	Little or no initial data exploration performed. Little or no relevancy. No or inappropriate use of Python.
25%	<b>3. Data Pre-processing</b>	Exceptionally thorough and extensive consideration of data quality issues for pre-processing. Appropriate approaches adopted with exceptionally clear understanding. Excellent use of Python.	Thorough consideration of data quality issues for pre-processing. Appropriate approaches adopted with clear understanding. Excellent use of Python.	Good consideration of data quality issues for pre-processing. Appropriate approaches adopted with clear understanding and every aspect covered. Good and flexible use of Python.	Reasonable consideration of data quality issues for pre-processing. Appropriate approaches adopted with reasonable understanding and most of the main issues covered. Good use of Python.	Limited consideration of data quality issues for pre-processing. Some appropriate approaches adopted with limited understanding and limited coverage. Limited use of Python.	Inadequate and/or inappropriate view of data quality issues. Inappropriate approaches adopted. Poor use of Python.	Little or no data quality issues considered. Inappropriate approaches adopted. No or inappropriate use of Python.
15%	<b>4. Modelling</b>	Appropriate algorithms employed with exceptionally clear understanding. Modelling with excellent working knowledge of Python	Appropriate algorithms employed with clear understanding. Modelling with excellent working knowledge of SAS Enterprise Miner.	Appropriate algorithms employed with clear understanding. Good and flexible use of Python.	Appropriate algorithms employed with reasonable understanding. Good use of Python.	Some appropriate algorithms employed with limited understanding. Limited use of Python.	Inappropriate and/or inadequate algorithms employed. Poor use of Python.	Little or no algorithms employed. Little or no use of Python.
20%	<b>5. Model Evaluation</b>	Exceptionally thorough and clear model interpretation and comparison with regards to business concerns. Excellent and meaningful models/patterns created.	Thorough and clear model interpretation and comparison with regards to business concerns. Excellent meaningful models/patterns created.	Clear model interpretation and comparison with regards to business concerns. Significantly meaningful models/patterns created.	Basic model interpretation and comparison with regards to business concerns. Reasonable models/patterns created.	Weak model interpretation and comparison with regards to business concerns. Very limited meaningfulness. Probably lack some clarity.	Poor model interpretation and comparison with regards to business concerns. No or little meaningful models/patterns provided.	Little or no model interpretation and comparison with regards to business concerns.
20%	<b>6. Report</b>	Exceptionally clear and concise summary of project findings. May raise questions for future research. Exceptional outstanding presentation. Clear structure and layout.	Very clear and concise summary of project findings. May raise questions for future research. Outstanding presentation. Clear structure and layout.	Clear and concise summary of project findings. Excellent presentation. Clear structure and layout.	Clear review and summary of project findings. Good presentation with proper structure and layout.	Adequate review of project findings. Probably lack of some clarity. Acceptable presentation.	Inadequate review of project findings. Lack of clarity and accuracy. Poor presentation.	Little or no review of project findings. Significantly Lack of clarity and accuracy. Very poor presentation.

## How to get help

We will discuss this Coursework Specification in class. However, if you have related questions, please contact me [name and email] as soon as possible.

## Resources

List resources such as background reading, templates, samples, tools, videos, links, etc

## Quality assurance of coursework specifications

Coursework specifications within CSI division go through internal (for new modules with 100% coursework also through external) moderation. This is to ensure high quality, consistency and appropriateness of the coursework as well as to share best practice within the CSI division.