

Coursework Specification

CW_Specification_CSI_6_DMA_25/26, Coursework 1

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.

Module Code	CSI_6_DMA
Module Title	Data Mining and Big Data Analytics
Lecturer	Daqing Chen; George Bamfo; Kamal Thapa
% of Module Mark	70%
Distributed	08/10/2025
Submission Method	Submit online via this Module's Moodle site
Submission Deadline	17:00, 10/12/2025
Release of Feedback & Marks	Feedback and provisional marks will be available in Grades on Moodle from 05/01/2026

Coursework Aim:

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

Coursework Details:

Type:	Project report
Overall View:	<p>This assignment is to be undertaken individually. You should plan your work carefully following the module's weekly teaching programme (See the Module Guide) and have a regular discussion with your tutor to address any questions and issues you may have during this project. You are expected to produce a written report for this data mining project.</p> <p>The assignment involves analysing a real-world dataset to identify its underlying patterns and trends by using appropriate data mining techniques and algorithms. These patterns and trends are intended</p>

	<p>to be used to address certain business concerns. A dataset will be assigned to you by your tutor.</p>
Tasks:	<p>You are required to undertake the following tasks in this project:</p> <ol style="list-style-type: none"> 1. <u>Business Understanding</u> <ol style="list-style-type: none"> 1.1. Download the dataset assigned to you from the module Moodle site along with the data description file (i.e., the metadata). 1.2. Read the data description to learn the nature of the dataset, such as what it is about, and what certain business context it is associated with, etc. 1.3. Examine the dataset within its business context and identify Three meaningful problems that potentially can be addressed by analysing the dataset. 1.4. Translate the business problems to appropriate data mining problems. 1.5. You may also refer to any articles in the literature related to the dataset under consideration. Possible sources for such articles include sciencedirect, researchgate, and IEEE Xplore. 2. <u>Data Understanding</u> <ol style="list-style-type: none"> 2.1. Determine the data type of each variable in the dataset and make correction to any incorrect data types identified. This must be completed first before proceeding any further. 2.2. Perform initial data exploration to get to know more about the dataset, such as the total number of instances, the number of attributes (variables), and explore the basic statistics of each attribute, including value range, average, standard deviation, skewness, kurtosis, and mode, etc. 2.3. Identify any data quality issues within the dataset, including missing values, outliers, extreme values, incomparable value ranges, and imbalanced classes, etc. 2.4. Determine if the dataset is appropriate to be used for addressing the business problems you have identified in Task 1. If not, re-do Task 1. 3. <u>Data Preparation</u> <ol style="list-style-type: none"> 3.1. Choose appropriate methods for data pre-processing, which includes changing data

	<p>types, dealing with missing values, tackling outliers, extreme values, and imbalanced classes, conducting data transformation and normalisation, and reducing dimensionality, etc., wherever appropriate.</p> <p>3.2. Identify correlations among certain variables using proper metrics, e.g., Pearson's correlation.</p> <p>3.3. Determine and discuss which, why, and how each attribute should (not) be used in your analysis.</p> <p>3.4. Divide the whole dataset into several subsets to be used for training, test and validation in predictive modelling.</p> <p>4. Modelling</p> <p>4.1. Use the pre-processed dataset to perform the data mining tasks you have identified in Task 1.</p> <p>4.2. Choose appropriate techniques and algorithms for your analysis: Choose either k-means clustering OR association rule analysis for descriptive modelling; Choose either decision tree OR regression for predictive modelling.</p> <p>4.3. Determine appropriate settings of the algorithms to be applied, e.g., how many clusters to use in k-means clustering.</p> <p>4.4. Re-do data preparation in Task 3 if needed.</p> <p>5. Evaluation</p> <p>5.1. Provide an explicit and concise interpretation and explanation of the descriptive model.</p> <p>5.2. If k-means clustering is chosen as your descriptive model, the model interpretation typically includes the statistics of the samples in each cluster in respect to the variables involved, how each cluster is featured, what common features are shared among the samples in each cluster, and how the common features vary comparatively from one cluster to another.</p> <p>5.3. If association rule analysis is chosen as your descriptive model, the interpretation typically includes what co-occurrence each association rule represents, and the meaningfulness and usefulness of each association as measured by support, confidence and lift.</p>
--	---

	<p>5.4. Evaluate the performance of the predictive model in terms of various measures applicable, such as accuracy, SSE (sum of squared errors), generalisation ability (overfitting), simplicity and cost, etc. For decision tree, the corresponding decision rules should be provided and explained along with a tree diagram.</p> <p>5.5. Discuss how the descriptive and predictive models created can be used to address the original business problems identified in Task 1.</p> <p>6. Report</p> <p>6.1. Summarise your main findings from the project.</p> <p>6.2. Don't repeat any basic concepts and methods which can be found from a relevant textbook and/or lectures, for example, what is <i>k</i>-means clustering, and what is a decision tree is, etc. Instead, you need to demonstrate your understanding of the concepts and methods and how you have applied them in your own analytical project.</p>
Word Count:	<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>In your report, do not include and explain any basic concepts of the subject; Instead, you should demonstrate your understanding of the concepts by applying them appropriately in your coursework.</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>
Presentation:	<ul style="list-style-type: none"> • Report must contain Title page, Table of Contents, Abstract, Conclusion, and References. • Work must be referenced, and a bibliography provided. • Work must be submitted as a Word document (.doc/docx) or a PDF. • Course work must be submitted using Arial font size 11 (or larger if you need to), with a minimum of 1.5 line spacing. • Your student number must appear at the front of

	the coursework. Your name must not be on your coursework.
Referencing:	Harvard Referencing should be used, see your Library Subject Guide for guides and tips on referencing.
Regulations:	<p>Make sure you understand the University Regulations on expected academic practice and academic misconduct. Note in particular:</p> <ul style="list-style-type: none"> ▪ 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. ▪ 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. <p>TURNITIN: 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%	1. Business Understanding	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%	2. Data Understanding	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%	3. Data Pre-processing	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 outstanding 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.	
20%	4. Modelling	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 Python.	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.	
15%	5. Model Evaluation	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%	6. Report	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.