

# UECM3993 GROUP ASSIGNMENT

---

COURSE CODE & COURSE TITLE: UECM3993 PREDICTIVE MODELLING  
COURSE: AM, AS, FM DEPARTMENT: DMAS

---

## Instructions

1. This is a group assignment with **four** to **seven** students including a **group leader** per group.
2. **Group leader** need to submit the following items to `liewhh@utar.edu.my`:
  - a list of members (with signatures)
  - group title/name (cannot be too bizarre or offensive)
  - the dataset of interest from the given listfor documentation before the start of assignment during class.
3. One submission per group. **Group leader** is responsible to submit the following documents through email to `liewhh@utar.edu.my`:
  - (a) “Group Name” Report.pdf ..... Wednesday of Week 11
  - (b) “Group Name” program code(s) ..... Wednesday of Week 11
4. **Deadline of submission for group assignment report and group programming code** is 4.00pm, 31 March 2021 (Wednesday of Week 11).
5. **Group Presentation** will be scheduled in week 11 to 13, date and time to be announced. Each presentation is limited to a maximum of 12 minutes (10 groups per lecture).
6. In the case of **late submission** for the report and program script, 10% of the maximum marks will be deducted if the work is up to one day late (24 hours) and additional 10% of the maximum marks for each of the subsequent days.
7. **Plagiarism is not allowed**. If the works are found to be plagiarised, no marks will be given and the incident will be reported to the university for further action.
8. The group assignment report **must** contain the **measurement of contribution of each member to the project** in ratio or percentage. The total of percentage will be 100%.

# Marks

- The group assignment report must have a section on **individual contributions** (in the first page or second page or the appendix).

- If the section **does not exist**, all individual member will receive

$$0.75 \times \text{teamwork marks}.$$

- If the section **exists**, each member will receive

$$0.7 \times \text{teamwork marks} + 0.3 \times \frac{\text{member contribution index}}{\text{max member contribution index}}$$

- Each member will receive equal marks for the group programming code well unless the **group leader** wants to have a different weights for the group members.
- Each member will receive equal marks for the group oral presentation with extra marks for members who present really well unless the **group leader** wants to have a different weights for the group members.
- A group leader can be **re-elected** if more than half of the members are not happy with the group leader one week before the submission of the assignment.

# Group Assignment Report (24%)

1. Pick a dataset from the following list and perform **unsupervised** and **supervised** learning on them:

- Academic Performance Evolution for Engineering Students Data Set (challenging) (<https://data.mendeley.com/datasets/83tcx8psxv/1>, Objective: Predict the \_PRO results)
- Bank Marketing Data Set (<https://archive.ics.uci.edu/ml/datasets/Bank+Marketing>)
- Car Evaluation Data Set (<https://archive.ics.uci.edu/ml/datasets/Car+Evaluation>)
- Census Income Data Set (<https://archive.ics.uci.edu/ml/datasets/Census+Income>)
- Credit Approval Data Set (<https://archive.ics.uci.edu/ml/datasets/Credit+Approval>)
- Early stage diabetes risk prediction Data Set (<https://archive.ics.uci.edu/ml/datasets/Early+stage+diabetes+risk+prediction+dataset.>)
- Facebook metrics Data Set (<https://archive.ics.uci.edu/ml/datasets/Facebook+metrics>)
- Heart Disease Data Set (<https://archive.ics.uci.edu/ml/datasets/Heart+Disease>)
- Human Activity Recognition Using Smartphones Data Set (<https://archive.ics.uci.edu/ml/datasets/Human+Activity+Recognition+Using+Smartphones>)
- Internet Advertisements Data Set (challenging) (<https://archive.ics.uci.edu/ml/datasets/Internet+Advertisements>)
- Online Shoppers Purchasing Intention Dataset (<https://archive.ics.uci.edu/ml/datasets/Online+Shoppers+Purchasing+Intention+Dataset>)
- Polish Companies Bankruptcy Data Set (<http://archive.ics.uci.edu/ml/datasets/Polish+companies+bankruptcy+data>)
- SMS Spam Collection (<https://archive.ics.uci.edu/ml/datasets/SMS+Spam+Collection>)
- Statlog (German Credit Data) Data Set ([https://archive.ics.uci.edu/ml/datasets/statlog+\(german+credit+data\)](https://archive.ics.uci.edu/ml/datasets/statlog+(german+credit+data)))
- Website Phishing Data Set (<https://archive.ics.uci.edu/ml/datasets/Website+Phishing>)
- Difficult: selfBACK Data Set (<https://archive.ics.uci.edu/ml/datasets/selfBACK>)
- Difficult: Facebook Comment Volume Dataset (<https://archive.ics.uci.edu/ml/datasets/Facebook+Comment+Volume+Dataset>)
- Difficult: Dota2 Games Results Data Set (<https://archive.ics.uci.edu/ml/datasets/Dota2+Games+Results>)

- Others data set from <https://archive.ics.uci.edu/ml/datasets.php> will need lecturer approval.
2. Marks for popular data will be lower: Dataset with only one group taking will get 1.5 marks, with two groups taking will get 1.2 marks, with three groups will get 0.8 marks, with four groups taking will get 0.3 marks and with more than four groups will receive zero marks.
  3. By using appropriate statistical software framework (R, Python with Scikit-learn, WEKA, C++, etc.), build models with the different statistical learning approaches (both unsupervised and supervised learning methods) which are covered in this course (and those which are not covered, but descriptions and documentations are required for methods not introduced in lecture with good references), find the “best” model for the data set selected and make sure that the objectives are met.

## Group Programming Code (12%)

1. Write a programming code (or nicely structured programming codes) with an appropriate use of libraries which analyse the raw dataset which is picked in the group assignment report and works in a data science pipeline.
2. Marks will be **deducted** if the programming code is in notebook and/or markdown formats which are not directly executable.
3. The programming code can only use free and legal software. The default is R (and Python). The group who try other free and legal open source software (such as Java, C++) which are cross-platform and does not have too much dependencies, i.e. the program can run on Microsoft Windows (of various versions), GNU/Linux platform, MacOS/X, etc., will receive extra marks.
4. The programming code(s) need to demonstrate the appropriate use of **supervised** and **unsupervised** learning with the free and legal statistical software tool.

## Group Oral Presentation (12%)

1. Prepare presentation slides which summarises the group assignment report and possible future improvements.
2. An oral presentation which involves every member or a representative member is allowed.
3. The oral presentation should cover the following aspects:
  - A good description of the problem and a systematic use of unsupervised and supervised learning methods to discover important information from the dataset.
  - A good illustration of results and conclusions.
  - Well-timed and interesting presentation (heavy marks may be deducted for over-time presentation)