

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, 6 April 2022 (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 (8 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

- Marks will be equally distributed by default. If the group assignment report has a section on **individual contributions** (in the first page or second page or the appendix), each member will receive

$$\text{teamwork marks} \times \left(1 - 0.4 \times \frac{\text{max IC} - \text{IC}}{\text{max IC}}\right)$$

where **IC** = individual contribution. For example, (Note: the same contribution can be applied to programming code.)

- A group with 7 members with contributions (30%, 30%, 30%, 2.5%, 2.5%, 2.5%, 2.5%) and the report is 10 out of 12
 - * 3 members will get $10 \times \left(1 - 0.4 \times \frac{30-30}{30}\right) = 10$ marks
 - * 4 members will get $10 \times \left(1 - 0.4 \times \frac{30-2.5}{30}\right) = 6.33$ marks
- A group with 4 members with contributions (A:10%, B:20%, C:30%, D:40%) and the report is 10 out of 12:
 - * member A gets $10 \times \left(1 - 0.4 \times \frac{40-40}{40}\right) = 10$ marks
 - * member B gets $10 \times \left(1 - 0.4 \times \frac{40-30}{40}\right) = 9$ marks
 - * member C gets $10 \times \left(1 - 0.4 \times \frac{40-20}{40}\right) = 8$ marks
 - * member D gets $10 \times \left(1 - 0.4 \times \frac{40-20}{40}\right) = 7$ marks
- 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 (18%)

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)
 - Higher Education Students Performance Evaluation Data Set (<https://archive.ics.uci.edu/ml/datasets/Higher+Education+Students+Performance+Evaluation+Dataset>)
 - Non Verbal Tourists Data Set (<https://archive.ics.uci.edu/ml/datasets/Non+verbal+tourists+data>)
 - Raisin Data Set (<https://dergipark.org.tr/tr/download/article-file/1227592>)
 - Dry Bean Data Set (<https://archive.ics.uci.edu/ml/datasets/Dry+Bean+Dataset>)
 - Mushroom Data Set (<https://archive.ics.uci.edu/ml/datasets/mushroom>)
 - Russian Corpus of Biographical Texts Data Set (<https://archive.ics.uci.edu/ml/datasets/Russian+Corpus+of+Biographical+Texts>)
 - Shoulder Implant X-Ray Manufacturer Classification Data Set (<https://archive.ics.uci.edu/ml/datasets/Shoulder+Implant+X-Ray+Manufacturer+Classification>)
 - Exasens Data Set (<https://archive.ics.uci.edu/ml/datasets/Exasens>)
 - Clickstream Data for Online Shopping (<https://archive.ics.uci.edu/ml/datasets/clickstream+data+for+online+shopping>)
 - Polish Companies Bankruptcy Data Set (<http://archive.ics.uci.edu/ml/datasets/Polish+companies+bankruptcy+data>)
 - 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>)
 - Others data set from Kaggle, etc. which are relevant to AS, AM or FM programme which are not too simple (or too difficult) can be used but needs lecturer's approval.
2. $\max\{0, 0.5 - 0.1n\}$ mark will be awarded to the data with are investigated by n groups.
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 (10%)

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 and/or non-text 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 (10%)

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.
 - **Explain the algorithm** behind the best model with respect to following aspects:
 - The mathematical/statistical idea behind best supervised learning model for the problem.
 - Explain how the model would be updated if new data comes in.
 - Well-timed and interesting presentation (heavy marks may be deducted for over-time presentation)