



School of Graduate  
and Professional  
Education



<b>Module</b>	<b>ITC 6103 – APPLIED MACHINE LEARNING</b>		
<b>Term</b>	<b>WINTER SEMESTER 2024</b>		
<b>Assessment</b>	<b>PROJECT</b>	<b>Weight</b>	<b>50%</b>
<b>Duration</b>			
<b>Deliverables</b>	<ol style="list-style-type: none"> <li>1. Report in Turnitin</li> <li>2. Code in Blackboard</li> <li>3. An oral examination/ presentation of your work</li> <li>4. Code in GitHub</li> </ol>		
<b>Method of Submission</b>	<i>TurinitIn, Blackboard, GitHub</i>		
<b>Deadline:</b>	<i>13<sup>th</sup> Week</i> <i>US Grading scale</i>		

---

*The rules of academic ethics apply when taking this assessment, including the requirement that you produce work **without improper or unauthorized assistance** from anyone.*

---

## General Instructions

Your project involves a series of experiments, observations coming out of the experiments, and drawing conclusions. Essentially you will collect data (or it will be provided by the instructor), then a programming language will be used (you are encouraged to use python) along with the appropriate libraries to process the data. Tables, diagrams, and data visualizations are essential for presenting your findings.

**Deliverables:** a) code in blackboard in Python, along with instructions for running it b) a report of 3000±500 words that will present your findings, and will be submitted at Turnit-in. The report must be self-contained, that is all experiments performed and all conclusions should be reported. If you need to exceed the word limit, use an appendix. c) an oral presentation d) code in GitHub

**Team size:** three persons

		Person being rated		
		Person-1	Person-2	Person-3
Person doing the rating	Person-1	1.25	1	0.75
	Person-2	1.10	1.10	0.80
	Person-3	1	1	1
Average Rate		1.12	1.03	0.85
Individual score (project grade: 80%)		89.6	82.4	68

**Grading:** peer-review

Teams consist of two or three persons. Group members will be asked to rate the relative contribution of themselves and the other group member(s). The ratings provided by each member must add up to the number of persons the group consists of (see example above). These ratings will be taken into account in the final grading of the project for each individual.

**Example:** In a group consisting of three members, each member provides a rating of all group members. As this is a three-member group, the ratings provided by each member add up to 3.00. A rating of 1.00 means that the person in question did exactly as much as expected of him/her. A rating that is less than 1.00 means that the person in question did less than expected, whereas a rating that is greater than 1.00 means that this person's contribution was greater than expected.

---

*The rules of academic ethics apply when taking this assessment, including the requirement that you produce work **without improper or unauthorized assistance** from anyone.*

---

## General Instructions

**Plagiarism:** The act of using another's words, ideas or organizational patterns without crediting or acknowledging the source. This includes work generated or modified by AI. In particular use of generative AI that violates the instructor's articulated policy, or using it to complete an assessment (e.g. project, midterm) in a way not explicitly permitted by the instructor will be considered a breach of academic integrity.

## Specific instructions

To carry-out the classification, clustering, and regression task you may need to consider the following steps:

- a. Data description & Visualization that aids the comprehension of the problem.
- b. Data pre-processing.
- c. Data/feature selection/evaluation.
- d. Decide how to split the data between training and data set. (If not stipulated by the instructions)
- e. Use multiple classifiers and evaluate the parameters of each classifier: Try at least the following: Decision Trees, one based on ensemble learning (especially consider the Random Forests) and Neural Networks.
- f. Use clustering algorithms, evaluate parameters. Try at least two (e.g.: k-means, and agglomerative (hierarchical) clustering).
- g. In regression: Try at least linear regression, and polynomial regression. Explore regularization.
- h. Evaluate
  - a. the performance of each classifier: at least provide F1 measure, precision, recall and ROC curves (if applicable) and AUC.
  - b. clusters based on criteria such as silhouette, and inertia.
  - c. regression based on criteria such as the R score and others.
- i. Observe findings and draw conclusions.
- j. Future work: Also include things you might try/consider in the future.

### 1. Clustering: Market Segmentation: Unsupervised learning (25%)

#### 2. Source data & description: US Census Data (1990) Data Set:

<https://archive.ics.uci.edu/ml/datasets/US+Census+Data+%281990%29>

3. Summarize the data set by discovering clusters, evaluate and characterizing them.
4. Is there a cluster of outliers (i.e. a cluster that differs significantly from other clusters)?

---

*The rules of academic ethics apply when taking this assessment, including the requirement that you produce work **without improper or unauthorized assistance** from anyone.*

---

## 5. Regression (20%)

It is up to you to choose a regression problem, but you should inform the instructor and get approval for it. Indicative data sources: Kaggle.com

, <https://www.analyticsvidhya.com/>, <https://github.com/awesomedata/awesome-public-datasets>, <https://www.openml.org/>.

The data set should present some challenges, e.g., size, missing values, or categorical features.

## 6. Predicting outcome (30%)

Consider the following data which is about predicting whether the user clicked on an item i.e. whether there is the **clickout item** in the user session.

For that you will use various features). The data file will upload in a link shared with you.

1. In particular, there are users active in sessions and you will need to predict the existence of a 'clickout item'. There are various features: *user\_id, session\_id, timestamp, step, action\_type, reference, device, current\_filters, action\_type impressions*.

<b>Clickout_item</b>	Is the info we try to <u>predict</u> in a session of a user. If there two clickout items in a session, we wish to predict the last one.
<b>User_id, session_id</b>	A certain user in a certain session performs some actions
<b>Steps</b>	Action steps for a user session.
<b>Action_type</b>	Type of action performed (e.g. search for poi, interaction image etc.)
<b>Platform</b>	Location of the platform, country
<b>City</b>	Location of the platform, city
<b>Device</b>	Access device
<b>Current_filters</b>	Possible filtering
<b>Impressions</b>	Items seen by the user (note: one of them maybe the clickout item). Do <b>not</b> use this field

2. Explainable AI (XAI): The explainability of the machine learning models has become very important. In this task you are required to do some research on the Shapley Additive Explanation (SHAP) or the Local Interpretable Model-Agnostic (LIME). Examine the results of explainability on the current data set on the best model found in step 1.

### 1. Report Quality (15%)

---

The quality of report is based on many factors including: organization of the material, presentation of data, experiments, models, evaluation, drawing conclusion using various aids such as tables, diagrammes, equations etc., and references if applicable.

### 2. Oral Presentation (10%)

---

During the presentation each group will present their work in a comprehensive manner and will be called to answer questions regarding their work.

## Grading scale: US System

	GP	Letter	US
Excellent	4.00	A	90+
Very good	3.70	A-	86-89
Very good	3.50	B+	81-85
Good	3.00	B	73-80
Satisfactory	2.50	C+	64-72
Satisfactory	2.00	C	51-63
Fail	0	F	<50

---

*The rules of academic ethics apply when taking this assessment, including the requirement that you produce work **without improper or unauthorized assistance** from anyone.*

---