**Assignment 3 (20%)**
**STQD6324 Data Management**
**SEMESTER 2 2024/2025**

Using the `u.user` file from the MovieLens 100k Dataset, which can be downloaded from https://grouplens.org/datasets/movielens/, write a Python script that functions as a wrapper to execute Cassandra Query Language (CQL) and Spark2 Structured Query Language (SQL) in order to answer the following questions. For each question, display only the top ten results:

  i)   Calculate the average rating for each movie.
  ii)  Identify the top ten movies with the highest average ratings.
  iii) Find the users who have rated at least 50 movies and identify their favourite movie genres.
  iv)  Find all the users who are less than 20 years old.
  v)   Find all the users whose occupation is "scientist" and whose age is between 30 and 40 years old.

Your python script should include the following elements:
  1. Python libraries used to execute Spark2 and Cassandra sessions.
  2. Functions to parse the `u.user` file into HDFS.
  3. Functions to load, read, and create Resilient Distributed Dataset (RDD) objects.
  4. Functions to convert the RDD objects into DataFrames.
  5. Functions to write the DataFrame into the Keyspace database created in Cassandra.
  6. Functions to read the table back from Cassandra into a new DataFrame.

**Optional: You may also attempt the above questions using HBase and MongoDB.**

The deadline for submitting your script is **2025-06-19**. Please share your Jupyter Notebook with markdown via **GitHub**.

| Criteria | Marks | | |
| --- | --- | --- | --- |
| **Reproducibility** | 3<br>The notebook is<br>100% reproducible | 2<br>The notebook is<br>reproducible with a few missing<br>steps | 1<br>The notebook is<br>not reproducible |
| **Interpretation** | **15**<br>The interpretation of the findings is clear,<br>easily understandable, and logical | **10**<br>The interpretation of the findings<br>is mostly clear and<br>understandable, with minor<br>areas needing clarification | **5**<br>The interpretation of<br>the findings is unclear<br>and difficult to<br>understand, lacking<br>logical coherence |
| **Overall<br>GitHub<br>presentation** | 2<br>The overall GitHub is<br>i. properly structured,<br>ii.each section neatly organized,<br>iii. easy to follow | 1<br>Part of the GitHub is<br>i. properly structured,<br>ii.neatly organized,<br>iii. easy to follow | 0<br>The GitHub is<br>i. poorly structured,<br>ii. each section is not<br>organized,<br>iii. hard to follow |