

**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 **2025-06-28**. Please share your Jupyter Notebook with markdown via **GitHub**.

Criteria	Marks		
<b>Reproducibility</b>	<p>3</p> <p>The notebook is 100% reproducible</p>	<p>2</p> <p>The notebook is reproducible with a few missing steps</p>	<p>1</p> <p>The notebook is not reproducible</p>
<b>Interpretation</b>	<p>15</p> <p>The interpretation of the findings is clear, easily understandable, and logical</p>	<p>10</p> <p>The interpretation of the findings is mostly clear and understandable, with minor areas needing clarification</p>	<p>5</p> <p>The interpretation of the findings is unclear and difficult to understand, lacking logical coherence</p>
<b>Overall GitHub presentation</b>	<p>2</p> <p>The overall GitHub is</p> <ul style="list-style-type: none"> <li>i. properly structured,</li> <li>ii. each section neatly organized,</li> <li>iii. easy to follow</li> </ul>	<p>1</p> <p>Part of the GitHub is</p> <ul style="list-style-type: none"> <li>i. properly structured,</li> <li>ii. neatly organized,</li> <li>iii. easy to follow</li> </ul>	<p>0</p> <p>The GitHub is</p> <ul style="list-style-type: none"> <li>i. poorly structured,</li> <li>ii. each section is not organized,</li> <li>iii. hard to follow</li> </ul>