# Question-2

We have created schema for the tables in Q1 along with basic cleaning, and now we want to create dimensional and fact tables for the dimensional modelling. To do so, we have to define their schemas as done in Q1 along with the fact table.

Since, it is a dimensional modelling, we have to pre-process the data and perform joins to form the fact table. Thus, some additional pre-processing will be required, specially on the subject names as all the three files have different structures of defining them altogether. Typically, fact tables are a formed as a result of join on student id and course/subject name. Thus, the course/subject name needs to be pre-processed in the hiveql, so it has similar contents which essentially belong to the same course. The data is from various sources like erp, codetantra and/or LMS, thus creating different values for the same subjects.

The typical pre-processing steps that we had one are in pre-processing.hql. The details of these pre-processing along with reasoning are as follows:

1. **Standardization of Text Format**

- **Convert to Lowercase:**

   - Using functions like LOWER() to convert all subject/course names to lowercase.
   - Reason: Ensures that differences in capitalization (e.g., "Maths" vs. "maths") do not create duplicate keys.

- **Trim Whitespaces:**

   - Apply the TRIM() function to remove leading and trailing spaces from text fields, especially trimming around delimiters like **/** and **-**, as evident in enrollment and grade data course fields.
   - Reason: Removes accidental spaces that could lead to mismatches during joins.

- **Uniform Delimiter Replacement:**

   - Use REGEXP_REPLACE() to standardize delimiters (like replacing hyphens, slashes, or multiple spaces with a single delimiter /). This is done to separate course code with course name.
   - Reason: Multiple representations (e.g., "CSE-101", "CSE/101", "CSE 101") get unified to a single format.

2. **Issues with attendance data:**

- **Uneven course_name in attendance data:**
   - Courses like "T1-24-25-AMS 211-Mathematics-3" are there in those fields, which should be ideally be "AMS 211-Mathematics-3" to maintain homogenity with other tables.
   - There are multiple rows which has email as **vishnu.raj@iiitb.org**. Those columns are essentially faculty meetings, and those rooms are removed and added to error_logs table, since they are erroneous values.
   - Some course names do not have any course code, and are essentially random staff/board meeting like **Audio testing Meeting by Prof Chandrashekar Ramanathan**. Those rows are removed from the table and added into error_logs table, since they are erroneous values.
- **Courses specifying batches:**
   - For courses with regard to first years, in some places they have mentioned batches they are teaching like **T1-24-25-GNL 101-English(BT1-IMT1-CSE)**. So, I have removed the contents of

the brackets except the ones which are programming courses like **T1-24-25-EGC 111-Programming 1A (C Programming)(BT1-IMT1)**.

3. **Final Course Details:**

- Now all the dimension tables have courses like **AMS 211/Mathematics-3**, meaning / is seperating the course code and course name.
  - Now, we could not merge directly with course codes since many rows are those of **programming and labs which have same course code**, but should have seperate grading and attendace records. Thus, standardisation of data across all the tables was required.

The hql queries for pre-processing is in **pre-processing.hql**.

Some images regarding sql queries done for pre-processing and data analytics are as follows:-

```
0: jdbc:hive2://localhost:10000/> SELECT DISTINCT `exam_result` FROM grade_roster;
INFO  : Compiling command(queryId=hive_20250414175135_8f56ea9a-7aa0-46d1-a62b-30b6820f4fe8): SELECT DISTINCT `exam_result` FROM grade_roster
INFO  : No Stats for student_data@grade_roster, Columns: exam_result
INFO  : Semantic Analysis Completed (retrial = false)
INFO  : Created Hive schema: Schema(fieldSchemas:[FieldSchema(name:exam_result, type:string, comment:null)], properties:null)
INFO  : Completed compiling command(queryId=hive_20250414175135_8f56ea9a-7aa0-46d1-a62b-30b6820f4fe8); Time taken: 0.112 seconds
INFO  : Concurrency mode is disabled, not creating a lock manager
INFO  : Executing command(queryId=hive_20250414175135_8f56ea9a-7aa0-46d1-a62b-30b6820f4fe8): SELECT DISTINCT `exam_result` FROM grade_roster
INFO  : Query ID = hive_20250414175135_8f56ea9a-7aa0-46d1-a62b-30b6820f4fe8
INFO  : Total jobs = 1
INFO  : Launching Job 1 out of 1
INFO  : Starting task [Stage-1:MAPRED] in serial mode
INFO  : Subscribed to counters: [] for queryId: hive_20250414175135_8f56ea9a-7aa0-46d1-a62b-30b6820f4fe8
INFO  : Tez session hasn't been created yet. Opening session
INFO  : Dag name: SELECT DISTINCT `exam......FROM grade_roster (Stage-1)
INFO  : HS2 Host: [ecb0cf9a7ce1], Query ID: [hive_20250414175135_8f56ea9a-7aa0-46d1-a62b-30b6820f4fe8], Dag ID: [dag_1744653095373_0001_1], DAG Session ID: [application_1744653095373_0001]
INFO  : Status: Running (Executing on YARN cluster with App id application_1744653095373_0001)

INFO  : Completed executing command(queryId=hive_20250414175135_8f56ea9a-7aa0-46d1-a62b-30b6820f4fe8); Time taken: 0.673 seconds
+--------------+
| exam_result  |
+--------------+
|              |
| Pass         |
| NULL         |
+--------------+
```

```
INFO  : Executing command(queryId=hive_20250414175713_cb5a9dd1-695e-41a2-818b-896f61aefd89): SELECT *
FROM attendance_data
LIMIT 5
INFO  : Completed executing command(queryId=hive_20250414175713_cb5a9dd1-695e-41a2-818b-896f61aefd89); Time taken: 0.0 seconds
+----------------------------------------+------------------------------------------------+----------------------------------------------------------+--------------------------+
|            attendance_data.course       |             attendance_data.instructor         |           attendance_data.name          |        attendance_data.email_id       |
| attendance_data.member_id | attendance_data.number_of_classes_attended | attendance_data.number_of_classes_absent | attendance_data.average_attendance_percent |
+----------------------------------------+------------------------------------------------+----------------------------------------------------------+--------------------------+
| T1-24-25-EGC 223 -Computer Architecture - Memory | nanditha.rao@iiitb.ac.in | 46290f2925cd1c7f330d5ed482bf9bbc7089ad5f7dba280cea6fadc02cd27a15 | 831f9b7f23152de96c2e022ef2299fbd8fbd0972e9a16f98d1bcb7c09d70b82a | 68f2511222cbc32ad56175871e928fcadcc965eb7cb49e8648b14796b7b53f8c | 7 | 5 | 58.3 |
| T1-24-25-EGC 113-Signals and Systems | jbapat@iiitb.ac.in, vinod.reddy@iiitb.ac.in | 46290f2925cd1c7f330d5ed482bf9bbc7089ad5f7dba280cea6fadc02cd27a15 | 831f9b7f23152de96c2e022ef2299fbd8fbd0972e9a16f98d1bcb7c09d70b82a | 68f2511222cbc32ad56175871e928fcadcc965eb7cb49e8648b14796b7b53f8c | 18 | 9 | 66.7 |
| T1-24-25-EGC 211-Programming 2A (C++ Programming) | ajay.bakre@iiitb.ac.in | 46290f2925cd1c7f330d5ed482bf9bbc7089ad5f7dba280cea6fadc02cd27a15 | 831f9b7f23152de96c2e022ef2299fbd8fbd0972e9a16f98d1bcb7c09d70b82a | 68f2511222cbc32ad56175871e928fcadcc965eb7cb49e8648b14796b7b53f8c | 7 | 10 | 41.2 |
| T1-24-25-EGC 212-Programming 2B (Java Programming) | vivek.yadav@iiitb.ac.in | 46290f2925cd1c7f330d5ed482bf9bbc7089ad5f7dba280cea6fadc02cd27a15 | 831f9b7f23152de96c2e022ef2299fbd8fbd0972e9a16f98d1bcb7c09d70b82a | 68f2511222cbc32ad56175871e928fcadcc965eb7cb49e8648b14796b7b53f8c | 3 | 1 | 75.0 |
| T1-24-25-AMS 203P-Physics-Lab | malapaka@iiitb.ac.in, bashok@iiitb.ac.in | 46290f2925cd1c7f330d5ed482bf9bbc7089ad5f7dba280cea6fadc02cd27a15 | 831f9b7f23152de96c2e022ef2299fbd8fbd0972e9a16f98d1bcb7c09d70b82a | 68f2511222cbc32ad56175871e928fcadcc965eb7cb49e8648b14796b7b53f8c | 0 | 1 | 0.0 |
+----------------------------------------+------------------------------------------------+----------------------------------------------------------+--------------------------+
5 rows selected (0.1 seconds)
0: jdbc:hive2://localhost:10000/>
```

```
INFO  : Executing command(queryId=hive_20250414180416_04b03066-e10f-4758-8230-b396cdf8c2c4): select * from grade_roster limit 5
INFO  : Completed executing command(queryId=hive_20250414180416_04b03066-e10f-4758-8230-b396cdf8c2c4); Time taken: 0.0 seconds
+----------------------------------------+------------------------------------------------+------------------------------+------------------------+
|     grade_roster.academy_location      |            grade_roster.student_id             |  grade_roster.student_status |   grade_roster.admission_id   | grade_roster.admission_status | grade_roster.student_name | grade_roster.program_name | grade_roster.batch | grade_roster.period | grade_roster.subject_code_name | grade_roster.section | grade_roster.faculty_name | grade_roster.course_credit | grade_roster.obtained_marks_grade | grade_roster.out_of_marks_grade | grade_roster.exam_result |
+----------------------------------------+------------------------------------------------+------------------------------+------------------------+
| International Institute of Information Technology Bangalore | f2c567e727dd3730f75b90f59460f6ab50c975cb8ea0ea653403f783df0e67df | Active | e25f6ad16689e911e9387afcf722692df48871cba02360a18ccc9a539f2c0ae0 | Active | 3df44cae9ae9235bb8b98503f097760aaf49fc18f2518f4a43c82def66918471 | Master of Science by Research-Part time | 2023 | Term I [2024-25] | IT 989/8 | T124-NR-989-8 | Nanditha Rao | 8 | 5 | A | Pass |
| International Institute of Information Technology Bangalore | 5d7799b46ed3c5e2a0d45f6444cd91efda0c9d4c89ba7b2ca | 5d7799b46ed3c5e2a0d45f6444cd91efda0c9d4c89ba7b2cac9d8e65f5acdd94 | Active | 12324eb9355064d6e3971dd3c752137333610484eb4787045d9a563be2ee45a | Master of Science By Research | MS Batch of 2024 | Term I [2024-25] | IT 989/8 | T124-RB-989-8 | Raghuram Bharadwaj | 8 | 8 | A | Pass |
| International Institute of Information Technology Bangalore | cce428fe0b9d509ceb3b52094ad39a03ed4a4920287bde92ad1df2c28a0c9400 | Active | cce428fe0b9d509ceb3b52094ad39a03ed4a4920287bde92ad1df2c28a0c9400 | Active | 4f2a9af570cff2f24573a7157dec547fe312b28ab8ad4dfad315ea3d82626d17 | Master of Science By Research | MS Batch of 2024 | Term I [2024-25] | IT 989/8 | T124-RB-989-8 | Raghuram Bharadwaj | 8 | 8 | A | Pass |
| International Institute of Information Technology Bangalore | ffba274d8a68b64e86980a5d807a0057faa389d2c7a585742 | ffba274d8a68b64e86980a5d807a0057faa389d2c7a585742 | Active | 4d47dc960e8c434 | 67530dd8202b82b761bc3be2325e34b52ca009c0f49f76d3fa6426cbbac4b868 | Master of Science By Research | MS Batch of 2024 | Term I [2024-25] | IT 989/8 | T124-RB-989-8 | Raghuram Bharadwaj | 8 | A | A | Pass |
| International Institute of Information Technology Bangalore | 07f606352fdbbf70fce33c81fcdb781f8b4dd09d932a078cd9b6b80b6257884a | Active | 63fc02dd5181ffb81fe8b5d1287f7fcf1aadd98d1f3839fff19973029750d55a | Active | cd83b07a473292be024a9e5305165c334daa37cae1da4a013f9f57a225e85cc7 | Doctor of Philosophy-Part time | PhD 2016-17 | Term I [2024-25] | IT 999/8 | T122-TKS-999-8 | T K Srikanth | 8 | 5 | A | Pass |
+----------------------------------------+------------------------------------------------+------------------------------+------------------------+
5 rows selected (0.087 seconds)
0: jdbc:hive2://localhost:10000/>
```

/

```
INFO  : Compiling command(queryId=hive_20250414181729_6604b463-7cff-4a92-b5af-62508cb5f1c5): CREATE TABLE merged_table (
serial_no INT,
course_type STRING,
student_id STRING,
student_name STRING,
program STRING,
batch STRING,
period STRING,
enrollment_date STRING,
primary_faculty STRING,
subject_code_name_enrollment STRING,
section_enrollment STRING,
academy_location STRING,
student_status STRING,
admission_id STRING,
admission_status STRING,
student_name_grade STRING,
program_name STRING,
batch_grade STRING,
period_grade STRING,
subject_code_name_grade STRING,
section_grade STRING,
faculty_name STRING,
course_credit INT,
obtained_marks_grade STRING,
out_of_marks_grade STRING,
exam_result STRING
)
ROW FORMAT DELIMITED
FIELDS TERMINATED BY ','
STORED AS TEXTFILE
```

```
0: jdbc:hive2://localhost:10000/> INSERT OVERWRITE TABLE merged_table
. . . . . . . . . . . . . . . . > SELECT
. . . . . . . . . . . . . . . . >       e.serial_no,
. . . . . . . . . . . . . . . . >       e.course_type,
. . . . . . . . . . . . . . . . >       e.student_id,
. . . . . . . . . . . . . . . . >       e.student_name,
. . . . . . . . . . . . . . . . >       e.program,
. . . . . . . . . . . . . . . . >       e.batch,
. . . . . . . . . . . . . . . . >       e.period,
. . . . . . . . . . . . . . . . >       e.enrollment_date,
. . . . . . . . . . . . . . . . >       e.primary_faculty,
. . . . . . . . . . . . . . . . >       e.subject_code_name AS subject_code_name_enrollment,
. . . . . . . . . . . . . . . . >       e.section AS section_enrollment,
. . . . . . . . . . . . . . . . >       g.academy_location,
. . . . . . . . . . . . . . . . >       g.student_status,
. . . . . . . . . . . . . . . . >       g.admission_id,
. . . . . . . . . . . . . . . . >       g.admission_status,
. . . . . . . . . . . . . . . . >       g.student_name AS student_name_grade,
. . . . . . . . . . . . . . . . >       g.program_name,
. . . . . . . . . . . . . . . . >       g.batch AS batch_grade,
. . . . . . . . . . . . . . . . >       g.period AS period_grade,
. . . . . . . . . . . . . . . . >       g.subject_code_name AS subject_code_name_grade,
. . . . . . . . . . . . . . . . >       g.section AS section_grade,
. . . . . . . . . . . . . . . . >       g.faculty_name,
. . . . . . . . . . . . . . . . >       g.course_credit,
. . . . . . . . . . . . . . . . >       g.obtained_marks_grade,
. . . . . . . . . . . . . . . . >       g.out_of_marks_grade,
. . . . . . . . . . . . . . . . >       g.exam_result
. . . . . . . . . . . . . . . . > FROM
. . . . . . . . . . . . . . . . >       enrollment_data e
. . . . . . . . . . . . . . . . > LEFT JOIN
. . . . . . . . . . . . . . . . >       grade_roster g
. . . . . . . . . . . . . . . . >       ON e.student_id = g.student_id
. . . . . . . . . . . . . . . . >       AND e.student_name = g.student_name
. . . . . . . . . . . . . . . . >       AND e.subject_code_name = g.subject_code_name;
```

```
INFO  : Query ID = hive_20250414181736_4d9ba148-0f5f-49dd-9ab4-fa8f4229c41f
INFO  : Total jobs = 1
INFO  : Launching Job 1 out of 1
INFO  : Starting task [Stage-1:MAPRED] in serial mode
INFO  : Subscribed to counters: [] for queryId: hive_20250414181736_4d9ba148-0f5f-49dd-9ab4-fa8f4229c41f
INFO  : Session is already open
INFO  : Dag name: INSERT OVERWRITE TABL......subject_code_name (Stage-1)
INFO  : Setting tez.task.scale.memory.reserve-fraction to 0.30000001192092896
INFO  : HS2 Host: [ecb0cf9a7ce1], Query ID: [hive_20250414181736_4d9ba148-0f5f-49dd-9ab4-fa8f4229c41f], Dag ID: [dag_1744654494947_0001_2],
INFO  : Status: Running (Executing on YARN cluster with App id application_1744654494947_0001)

INFO  : Starting task [Stage-2:DEPENDENCY_COLLECTION] in serial mode
INFO  : Starting task [Stage-0:MOVE] in serial mode
INFO  : Loading data to table student_data.merged_table from file:/opt/hive/data/warehouse/student_data.db/merged_table/.hive-staging_hive_
INFO  : Starting task [Stage-3:STATS] in serial mode
INFO  : Executing stats task
-------------------------------------------------------------------------------
        VERTICES      MODE        STATUS   TOTAL   COMPLETED  RUNNING  PENDING  FAILED  KILLED
-------------------------------------------------------------------------------
Map 3 .......... container     SUCCEEDED      1          1        0        0       0       0
Map 1 .......... container     SUCCEEDED      1          1        0        0       0       0
Reducer 2 ...... container     SUCCEEDED      1          1        0        0       0       0
-------------------------------------------------------------------------------
VERTICES: 03/03  [==========================>>] 100%  ELAPSED TIME: 0.99 s
-------------------------------------------------------------------------------
INFO  : Table student_data.merged_table stats: [numFiles=1, numRows=3634, totalSize=1539904, rawDataSize=1536270, numFilesErasureCoded=0]
INFO  : Completed executing command(queryId=hive_20250414181736_4d9ba148-0f5f-49dd-9ab4-fa8f4229c41f); Time taken: 1.055 seconds
3,634 rows affected (1.342 seconds)
0: jdbc:hive2://localhost:10000/>
```

What we did was first write the python script for all the dimensional tables and pre-processed it such that on doing inner join, we will get maximum rows in the fact table. Now, the fact table has **2771 rows**, which would have been **less than 1000** without pre-processing. Then, we backtracked and form the hql queries and reported it in hql file.

The structure of fact tables is as follows:

```
CREATE TABLE IF NOT EXISTS fact_table (
    member_id STRING,
    course STRING,
    number_of_classes_attended INT,
    number_of_classes_absent INT,
    course_credit INT,
```

```
      average_attendance_percent FLOAT
  )
  ROW FORMAT SERDE 'org.apache.hadoop.hive.serde2.OpenCSVSerde'
  WITH SERDEPROPERTIES (
    "separatorChar" = ",",
    "quoteChar"     = "\""
  )
  STORED AS TEXTFILE
  TBLPROPERTIES ("skip.header.line.count" = "1");
```

The structure of all the dimension tables as defined in Q1 are as follows:-

```
  CREATE TABLE IF NOT EXISTS dim_enrollment_data (
    serial_no INT,
    course_type STRING,
    student_id STRING,
    student_name STRING,
    program STRING,
    batch STRING,
    period STRING,
    enrollment_date STRING,
    primary_faculty STRING,
    subject_code_name STRING,
    section STRING
  )
  ROW FORMAT SERDE 'org.apache.hadoop.hive.serde2.OpenCSVSerde'
  WITH SERDEPROPERTIES (
    "separatorChar" = ",",
    "quoteChar"     = "\""
  )
  STORED AS TEXTFILE
  TBLPROPERTIES ("skip.header.line.count"="1");


  CREATE TABLE IF NOT EXISTS dim_grade_roster (
      academy_location STRING,
      student_id STRING,
      student_status STRING,
      admission_id STRING,
      admission_status STRING,
      student_name STRING,
      program_name STRING,
      batch STRING,
      period STRING,
      section STRING,
      faculty_name STRING,
      course_credit INT,
      obtained_marks_grade STRING,
      out_of_marks_grade STRING,
      exam_result STRING,
      subject_code_name STRING
  )
```

```
   ROW FORMAT SERDE 'org.apache.hadoop.hive.serde2.OpenCSVSerde'
   WITH SERDEPROPERTIES (
       "separatorChar" = ",",
       "quoteChar"     = "\""
   )
   STORED AS TEXTFILE
   TBLPROPERTIES ("skip.header.line.count"="1");

   CREATE TABLE IF NOT EXISTS dim_attendance_data (
       course STRING,
       instructor STRING,
       name STRING,
       email_id STRING,
       member_id STRING,
       number_of_classes_attended INT,
       number_of_classes_absent INT,
       average_attendance_percent FLOAT
   )
   ROW FORMAT SERDE 'org.apache.hadoop.hive.serde2.OpenCSVSerde'
   WITH SERDEPROPERTIES (
       "separatorChar" = ",",
       "quoteChar"     = "\""
   )
   STORED AS TEXTFILE
   TBLPROPERTIES ("skip.header.line.count"="1");
```

Firstly, we mount the csv files into the docker image folder, so as to use it for populating tables with the data.

```
keshav-chandak@keshav-chandak-HP-Pavilion-Laptop-14-ec1xxx:~/Desktop/output Q2$ docker cp attendance.csv hive4:/tmp/dim_attendance.csv
Successfully copied 2.36MB to hive4:/tmp/dim_attendance.csv
keshav-chandak@keshav-chandak-HP-Pavilion-Laptop-14-ec1xxx:~/Desktop/output Q2$ docker cp enrollment.csv hive4:/tmp/dim_enrollment.csv
Successfully copied 868kB to hive4:/tmp/dim_enrollment.csv
keshav-chandak@keshav-chandak-HP-Pavilion-Laptop-14-ec1xxx:~/Desktop/output Q2$ docker cp grade.csv hive4:/tmp/dim_grade.csv
Successfully copied 1.8MB to hive4:/tmp/dim_grade.csv
keshav-chandak@keshav-chandak-HP-Pavilion-Laptop-14-ec1xxx:~/Desktop/output Q2$ docker cp fact_table_final1.csv hive4:/tmp/fact_table.csv
keshav-chandak@keshav-chandak-HP-Pavilion-Laptop-14-ec1xxx:~/Desktop/output Q2$ 
```

Then, we load the csv dataset into the above schema.

The code for loading it into hql table schemas is in load_queries.hql

The corresponding hql output after loading, and select statements are as follows:

```
+----------------------------------------------+------------------------------------------------------+------------------------
-----------------------------------+---------------------------------------+------------------------
---------------------------------------+
8,495 rows selected (2.601 seconds)
```

```
----------------------------------------+----------------------------+
 3,101 rows selected (0.313 seconds)
```

```
---+------------------------------+----------------------------------------------+
 4,477 rows selected (0.478 seconds)
```

Fater this is done, we try three HiveQl analytic queries. I have utilised these three queries since it covers the utility of all the numerical columns in the dimension and fact tables.

Before starting off, since we are utilising hive as a docker image due to various issues in the instllation as faced by many others, we are storing the tables everytime in our local system. So, first we load csv of dimensional tables and fact table onto the docker image: docker cp attendance.csv hive4:/tmp/dim_attendance.csv

docker cp enrollment.csv hive4:/tmp/dim_enrollment.csv docker cp grade.csv hive4:/tmp/dim_grade.csv

/

docker cp fact_table_final.csv hive4:/tmp/fact_table.csv

```
keshav-chandak@keshav-chandak-HP-Pavilion-Laptop-14-ec1xxx:~/Desktop/output Q2$ docker cp attendance.csv hive4:/tmp/dim_attendance.csv
Successfully copied 2.36MB to hive4:/tmp/dim_attendance.csv
keshav-chandak@keshav-chandak-HP-Pavilion-Laptop-14-ec1xxx:~/Desktop/output Q2$ docker cp enrollment.csv hive4:/tmp/dim_enrollment.csv
Successfully copied 868kB to hive4:/tmp/dim_enrollment.csv
keshav-chandak@keshav-chandak-HP-Pavilion-Laptop-14-ec1xxx:~/Desktop/output Q2$ docker cp grade.csv hive4:/tmp/dim_grade.csv
Successfully copied 1.8MB to hive4:/tmp/dim_grade.csv
keshav-chandak@keshav-chandak-HP-Pavilion-Laptop-14-ec1xxx:~/Desktop/output Q2$ docker cp fact_table_final1.csv hive4:/tmp/fact_table.csv
keshav-chandak@keshav-chandak-HP-Pavilion-Laptop-14-ec1xxx:~/Desktop/output Q2$ 
```

## Query-1

**Objective:**

To compute the CGPA (Cumulative Grade Point Average) for each student based on the grade obtained and course credits.

**Approach:**

- Join `dim_grade_roster` and `fact_table` on `student_id` and `subject_code_name`.
- Use a weighted sum of grade points (based on institutional grading system) multiplied by `course_credit`.
- Divide total weighted grade points by total credits to derive CGPA.
- Order results by CGPA and then by total credits in descending order.

**Query**

```sql
SELECT
  g.student_id,
  SUM(g.course_credit) AS total_credits_completed,
  SUM(CASE
        WHEN g.obtained_marks_grade = 'A'  THEN 4.0 * g.course_credit
        WHEN g.obtained_marks_grade = 'A-' THEN 3.7 * g.course_credit
        WHEN g.obtained_marks_grade = 'B+' THEN 3.4 * g.course_credit
        WHEN g.obtained_marks_grade = 'B'  THEN 3.0 * g.course_credit
        WHEN g.obtained_marks_grade = 'B-' THEN 2.7 * g.course_credit
        WHEN g.obtained_marks_grade = 'C+' THEN 2.4 * g.course_credit
        WHEN g.obtained_marks_grade = 'C'  THEN 2.0 * g.course_credit
        WHEN g.obtained_marks_grade = 'D'  THEN 1.7 * g.course_credit
        ELSE 0.0
      END) / SUM(g.course_credit) AS cgpa
FROM dim_grade_roster g
JOIN fact_table f
  ON g.student_id = f.member_id
  AND g.subject_code_name = f.course
GROUP BY g.student_id
ORDER BY cgpa DESC, total_credits_completed DESC;
```

**Use Case:**

This query is essential for academic performance analysis, ranking students, and eligibility for honors or scholarships.

```
----------------------------------------------------------------------------------
     VERTICES     MODE        STATUS  TOTAL  COMPLETED  RUNNING  PENDING  FAILED  KILLED
----------------------------------------------------------------------------------
Map 3 .......... container   SUCCEEDED    1       1        0        0       0      0
Map 1 .......... container   SUCCEEDED    1       1        0        0       0      0
Reducer 2 ...... container   SUCCEEDED    1       1        0        0       0      0
----------------------------------------------------------------------------------
VERTICES: 03/03  [==========================>>] 100%  ELAPSED TIME: 2.34 s
----------------------------------------------------------------------------------
INFO  : Completed executing command(queryId=hive_20250414133532_77c2051d-3fbd-4fb1-92b1-c8a381e49198); Time taken: 5.485 seconds
+--------------------------------------------+----------------------+--------------------+
|                g.student_id                | total_credits_completed |       cgpa       |
+--------------------------------------------+----------------------+--------------------+
| 006ebffbd115df9b6ef0e30a5cf33a86d6544a0bdb4b2e0c5f01addf199fbe8f | 28.0 | 1.95              |
| 01021eb63ad8ca36d35a6fd4ead1a931e4dc4b74999a5cf98c7900d8540c97ae | 8.0  | 1.475             |
| 01104f71b9089725f8209bb949fb92555b90730dd4213561908386f1f0269a2b | 22.0 | 3.0090909090909084|
| 0133dbf630dcec089bb08ca3c4ec094ef4d383b985452330649c99a8acd5001a | 28.0 | 2.2857142857142856|
| 01e748f6f48344ff2bf1f20e5eb76b7411c8751af41798ff01d97fddae5d4234 | 12.0 | 3.233333333333333 |
| 03c401666f88bd87df6663255493524ba394e8db25ba9af794c9febc0c03f12b | 26.0 | 3.423076923076923 |
| 03e8af13a98d6f1287619ac0890c632fa203419b6f65a005c6c9d2f8478fe282 | 26.0 | 2.8999999999999995|
| 03f205b589909f0ea18950c4fac7e7d125a61a992e33556e8a3a8b0615ab0ab4 | 32.0 | 2.21875           |
| 047236cffacc85ccec880c7b1b257e321af0ef1dd290899de7d6f9319decda76 | 32.0 | 2.4               |
| 075e7f21e42b4a5fb6e97df2bd17e65a0af0e5b11f547bfecc4ca690a2ece98e | 32.0 | 3.1125            |
| 075f4288380a972f084731c23f3ae382165107e4c5a2a2cd85363d3a96046fed | 24.0 | 1.4166666666666667|
| 076449087afdae0e4172c37b1c10b693248751418392ad649ef57a52ad6e0e14 | 26.0 | 2.2923076923076926|
| 0821a962c2726e5df442dc86f74a371ce338c2436dd2e566f85f07883c5271c2 | 28.0 | 3.5357142857142856|
| 086ffcfc64ba1b317ff114d2d3dc632675ae75ee82788a8fa0b31e6be050394d | 28.0 | 3.7642857142857147|
| 08aa713e1d2c465191d99525020cf07f773e107a506a44229e7ffb500efd98dd | 26.0 | 1.4615384615384615|
```

```
| f7b37b09dd10930d9a0132e26d2830ca8677ad11d0f666db6fa0724fe57a1fff | 22.0 | 2.0636363636363635 |
| f800cfdc8d739f2d384761d93f76d5f8d4d5c24f8b63f96556a754e6c1f86c8c | 24.0 | 2.4166666666666665 |
| f853b03aaf270f8f8b6cab1ac5003975ffc2e14ce0f8d696e0f90d5c7e80421b | 32.0 | 3.09375            |
| f9746d5926e1ef8be988f4a01b8897189476a4792deee63c7aa37e2d31b862a3 | 22.0 | 2.9272727272727277 |
| f9a66b0efea2de779b86a5f40937feb83c080449e91f30eb7454b32d2d7295b6 | 22.0 | 1.7363636363636366 |
| f9c3e40f66a95bff6864d2daff1a29d32b55d0034e5753ae9095585f0202314e | 36.0 | 2.1055555555555556 |
| fa30950bc068d2bff9c983cb0853be94e0f15ba6fca5468c567db2ca275a7275 | 32.0 | 2.09375            |
| fa97ea0f7b79d3347a03f5cdc5e96188d59f7e7098a0cec26b28d2f804fcf205 | 32.0 | 2.03125            |
| fb82641a70b62444754aaca4126cf6d6566970fe04c5746b7f97312613a2f7fa | 32.0 | 3.375              |
| fbd0443bf1e0d231601b6aff94a29877222aca65946506425863c35151df2084 | 22.0 | 1.8636363636363635 |
| fc1e3958bf58979da2cd0fd53a5a62ba037f7eb11aebe44e08b2ea5f37cc2ffb | 36.0 | 2.138888888888889  |
| fc43072bf0449e0f4f3743a9fb44d63507c0444bf6db7440443111fb0f406bce | 28.0 | 3.1999999999999997 |
| fc4535a76a801757ff741a0cf4f9aef52866e36e06aacc43239945bd0cca113c | 28.0 | 2.9499999999999997 |
| fc5f93239ec1b27fd8bf7174a1f68e953d57e0b86e3c910135d02658a01a26ed | 26.0 | 1.9                |
| fcfa55660b5d441de2ef2e9b0b95b18c33a3f4853acdd231fea1eddd58dcc1ee | 22.0 | 1.6363636363636365 |
| fcfbf656fb89ac105f2d0a8393c61f314a8449184a2f72349eef90b477c6c37b | 12.0 | 1.9833333333333334 |
| fd9709ae2b08802a0cfc32aa1971dd29c0de7c8b4be3cc07a1cb968fe2405ed5 | 28.0 | 1.3642857142857143 |
| fdb1bf0b3ff8d8048103388f108794de4164bbe8bdbf7d898a6036965cc2f292 | 28.0 | 2.9285714285714284 |
| fe6cacdcebbf5892a3583e6ec13530f2e6ea7c6c75a90fcced9a2645e7200033 | 28.0 | 2.8928571428571423 |
| fedafcd150b9a17932760554a0ec9208266957a49da49214f4f9c7e1776f340d | 22.0 | 2.8636363636363633 |
| ff6358e8fa8dce631d81990d463738796e3eb5cb545a29edad662cd92864cbfb | 8.0  | 0.25               |
| ffba274d8a68b64e86980a5d807a0057faa389d2c7a5857424d47dc960e8c434 | 12.0 | 2.4166666666666665 |
| ffd48b5414c5c285193e34544de015ed643829e5bf39c79b107a5c41aaa612dd | 28.0 | 2.857142857142857  |
| ffe3d002fbf6b6c4020303b73c54bcef8c8e9c4b5db7108ac2c8f9b206f0f177 | 26.0 | 2.4461538461538463 |
+--------------------------------------------+----------------------+--------------------+
524 rows selected (6.54 seconds)
```

**Time Elapsed: 6.54 seconds**

## Query-2

**Objective:**
To determine the number of students taught, average attendance, and maximum course credit for each faculty.

**Approach:**

- Join `dim_grade_roster` and `fact_table` on student and course.
- Filter for only those students who have passed (`exam_result = 'Pass'`).
- Aggregate data to:
    - Count distinct students per faculty.
    - Calculate average attendance using `average_attendance_percent`.
    - Determine the highest credit course taught by each faculty.

**Use Case:**
This helps analyze faculty engagement, workload distribution, and effectiveness in teaching based on student attendance and course difficulty.

**Query:**

```sql
SELECT
  g.faculty_name,
  COUNT(DISTINCT g.student_id) AS num_students,
  AVG(f.average_attendance_percent) AS avg_attendance,
  MAX(g.course_credit) AS max_course_credit
FROM fact_table f
JOIN dim_grade_roster g
  ON f.member_id = g.student_id
    AND f.course = g.subject_code_name
WHERE g.exam_result = 'Pass'
GROUP BY g.faculty_name;
```

```
+------------------------+---------------+---------------------+-------------------+
|     g.faculty_name     | num_students  |    avg_attendance   | max_course_credit |
+------------------------+---------------+---------------------+-------------------+
| Amit Chattopadhyay     | 159           | 84.39371069182388   | 4.0               |
| Ashish Choudhury       | 6             | 80.73333333333333   | 4.0               |
| Badrinath Ramamurthy   | 120           | 87.2225             | 2.0               |
| G Srinivasa Raghavan   | 4             | 88.675              | 4.0               |
| Jaya Sreevalsan Nair   | 1             | 70.8                | 4.0               |
| Jyotsna Bapat          | 2             | 97.2                | 4.0               |
| Karthikeyan Vaidyanathan | 1           | 85.7                | 4.0               |
| Kurian Polachan        | 91            | 86.91978021978026   | 4.0               |
| Manisha Kulkarni       | 119           | 76.56722689075629   | 4.0               |
| Meenakshi D Souza      | 3             | 86.26666666666667   | 4.0               |
| Nanditha Rao           | 42            | 66.87857142857142   | 4.0               |
| Pillalamarri Sridhar   | 160           | 80.71               | 4.0               |
| Preeti Mudliar         | 33            | 80.2                | 4.0               |
| Priyanka Das           | 6             | 77.18333333333335   | 4.0               |
| Priyanka Sharma        | 280           | 66.44857142857144   | 2.0               |
| Prof. Amrita Mishra    | 120           | 79.95333333333339   | 4.0               |
| S Malapaka             | 166           | 80.00903614457827   | 4.0               |
| Sachit Rao             | 150           | 74.71743119266057   | 4.0               |
| Sakshi Arora           | 30            | 73.76666666666667   | 4.0               |
| Srinath Srinivasa      | 3             | 88.90000000000002   | 4.0               |
| Srinivas Vivek         | 198           | 77.55353535353527   | 4.0               |
| Sujit Kumar Chakrabrati| 160           | 86.43624999999997   | 2.0               |
| Sushree Behera         | 4             | 81.825              | 4.0               |
| Thangaraju B           | 149           | 92.32364864864864   | 4.0               |
| Tulika Saha            | 120           | 73.93666666666667   | 2.0               |
| Uttam Kumar            | 2             | 28.0                | 4.0               |
| V Sridhar              | 313           | 83.2861271676299    | 4.0               |
| Vinod Reddy            | 5             | 67.03999999999999   | 4.0               |
| Vinu E V               | 59            | 87.05762711864405   | 4.0               |
| Viswanath G            | 145           | 85.38620689655166   | 4.0               |
+------------------------+---------------+---------------------+-------------------+
30 rows selected (0.932 seconds)
```

**Time Elapsed:0.912 seconds**

## Query-3

**Objective:**
To identify students who have an attendance percentage below 75% in any course.

**Approach:**

- Join `dim_grade_roster` and `fact_table` on `student_id` and `subject_code_name`.
- Calculate overall attendance percentage as (classes_attended / (attended + absent)) * 100:
- Filter (`HAVING`) to return only those records with less than 75% attendance.

**Query**:

```
SELECT
    g.student_id,
    g.subject_code_name AS course,
    SUM(f.number_of_classes_attended) AS total_classes_attended,
    SUM(f.number_of_classes_absent)  AS total_classes_absent,
    (SUM(f.number_of_classes_attended) * 100.0) /
(SUM(f.number_of_classes_attended) + SUM(f.number_of_classes_absent)) AS
overall_attendance_percentage
FROM fact_table f
INNER JOIN dim_grade_roster g
  ON f.member_id = g.student_id
  AND f.course = g.subject_code_name
GROUP BY
    g.student_id,
    g.subject_code_name
HAVING
    (SUM(f.number_of_classes_attended) * 100.0) /
(SUM(f.number_of_classes_attended) + SUM(f.number_of_classes_absent)) < 75;
```

**Use Case:**

Used for academic warnings, eligibility checks for exams, and enforcing minimum attendance policies.



**Time Elapsed:1.23 seconds**

Note: You might be seeing that I am using only two tables in the queries, but since the fact table contains all the numerical data regarding attendance, thus **dim_attendance** table is not used. Similarly enrollment data had no numerical values, thus it is not part of join, as there cannot be any analytical query possible.

## Error logs

The error_log.csv in the output folder of Q2 contains the inconsistent and erroneous data that we found out earlier. Since, the rest of the data was pre-processed and retained in the table, only erroneous values in the attendance table has been copied to the error_logs table.

```sql
INSERT INTO TABLE error_log
SELECT *
FROM a_data
WHERE (course NOT REGEXP '[0-9]' OR email_id = 'vishnu.raj@iiitb.org');

INSERT OVERWRITE DIRECTORY '/tmp/error_log_csv'
ROW FORMAT DELIMITED
FIELDS TERMINATED BY ','
SELECT * FROM error_log;

docker cp hive4:/tmp/error_log_csv/000000_0 ./error_log.csv
```