

Question-4: Pig Query

Note: To run a pig query and make a time comparison, we can simply put:

```
time pig -x local sample.pig
```

The scripts that were used for querying are their in the **Q4 directory**.

The output from the pig query is there in the output directory

Query 1: CGPA Calculation per Student

Objective:

Calculate the CGPA for each student along with their total credits completed using the institutional grading system.

Steps and Logic:

1. Loading Data:

- Load `enrollment.csv` into `enrollment_data` and drop the header row.
- Load `grade.csv` into `grade_roster` and filter out its header row.
- Load `attendance.csv` and `fact_table_final1.csv` similarly, ensuring that header rows are filtered out.

2. Joining Datasets:

- Join `grade_roster` and `fact_table` on matching student and course identifiers (i.e. `student_id` with `member_id` and `subject_code_name` with `course`).

3. Calculating Weighted Points:

- Use a `CASE` expression within a `FOREACH` to calculate the weighted points for each course based on the letter grade (e.g., 'A' as 4.0, 'A-' as 3.7, etc.) multiplied by the course credit.

4. Grouping and Aggregation:

- Group the resulting data by `student_id` to aggregate values.
- Compute the total credits by summing `course_credit` and the total weighted points.
- Calculate the CGPA as the ratio of total weighted points to total credits.

5. Ordering and Storing:

- Order the results by CGPA (in descending order) and by total credits completed.
- Dump and store the output using PigStorage into the `/output/Query-1` directory.

Pig Script Excerpt:

```
-- Join grade_roster and fact_table
joined_data = JOIN grade_roster BY (student_id, subject_code_name),
```

```
fact_table BY (member_id, course);

-- Calculate weighted points
cgpa_data = FOREACH joined_data GENERATE
  grade_roster::student_id AS student_id,
  grade_roster::course_credit AS course_credit,
  (CASE grade_roster::obtained_marks_grade
    WHEN 'A' THEN 4.0 * grade_roster::course_credit
    WHEN 'A-' THEN 3.7 * grade_roster::course_credit
    WHEN 'B+' THEN 3.4 * grade_roster::course_credit
    WHEN 'B' THEN 3.0 * grade_roster::course_credit
    WHEN 'B-' THEN 2.7 * grade_roster::course_credit
    WHEN 'C+' THEN 2.4 * grade_roster::course_credit
    WHEN 'C' THEN 2.0 * grade_roster::course_credit
    WHEN 'D' THEN 1.7 * grade_roster::course_credit
    ELSE 0.0
  END) AS weighted_points;

-- Group and compute totals and CGPA
grouped_data = GROUP cgpa_data BY student_id;
result = FOREACH grouped_data {
  total_credits = SUM(cgpa_data.course_credit);
  total_weighted_points = SUM(cgpa_data.weighted_points);
  cgpa = total_weighted_points / total_credits;
  GENERATE group AS student_id, total_credits AS total_credits_completed,
  cgpa AS cgpa;
}
ordered_result = ORDER result BY cgpa DESC, total_credits_completed DESC;

DUMP ordered_result;
STORE ordered_result INTO '/output/Query-1' USING PigStorage(',');
```

Image

```
(6451f1fdbfb51e85c217fc58523ac177fe9c1e59b0dd11cf65d47035fdb720b1,28,2.3785714285714286)
(b0026ddcb2635476e72f335b7ded6341ede34eb5f8f3e2ed34aa60062d3934fd,32,2.375)
(01021eb63ad8ca36d35a6fd4ead1a931e4dc4b74999a5cf98c7900d8540c97ae,8,2.35)
(5db0ea3e96305d5b2434e0a9e3657a76c67509c8327c7e47fd54d1b4b9063f31,26,2.3307692307692305)
(4ca9268ea5658127bf8512445be6922e2c2357c8c52b5c5a2c631cae6af0c5d,22,2.2909090909090906)
(86b2b4629113bb3a78373aa25f95a8dccc6e5676327b2c3e36109872ac9bacc2c,12,2.2666666666666667)
(447f6ae3c7fd293dabbad856074c77f5ac90b133b9b114cc8080e78770d60882,22,2.2545454545454544)
(ad841dec8c5b2f1952dc51dd68ad5f13a672c4ddeaf228935127ce06c875c20a,22,2.2363636363636363)
(c4ff6797d1fb4433553653d82b0289c32a3c7837c7012da48ecab2ac9c01ee3a,32,2.1875)
(8777036516faed4eabf100af059a4c3e157ceaea6bcdafa28191d593415f64204,28,2.157142857142857)
(46afcb932a74d834d21d0547dee5b8bfefa616f4049b3b77dd47c045c2101cc9,32,2.1)
(cae5018bebe2ccff689e80a84c44d7d3a6acad9949e72e364e03c5e6116b9bed,32,2.1)
(19c5c74b92596d69c552ea81c9aab4370c91a8753e1b64020a7cf16367ea3ea1,26,2.0615384615384613)
(b7f97c4154f2b34129c7f5192e1c40e1436ec1f742924fddb0c19252dc5a15bb,32,2.025)
(9886071454e243de3ff7fcd672bc754d453b29a5f08e6c6e0df0c5cf8b47f4362,4,2.0)
(1d26c5cdc02e2b8d9ea7983ca28faf91d58e2755852f9ff5d2321a18e21d3b49,4,2.0)
(1b5b60677927228f94b20a68dad069e43e687a6e8eabeb81cff935eeab4f67,28,1.807142857142857)
(8177fca161b90d83cee14b5e9162f828670d8035a199b48bb5110432b623e9a7,22,1.7363636363636366)
(d24d0df17e3bbea38e89f2eda4c03d9b30c756f4fa6aeba9ff7e3f8cb7e78bc4,26,1.6846153846153844)
(08aa713e1d2c465191d99525020cf07f773e107a506a44229e7fffb500ef9d8dd,26,1.623076923076923)
(2b6a85597239bfe683c37419733fb3a9db6d8c4abfd93fa6db9db6f3bca9d493,28,1.5999999999999999)
(99f22d27d8c07ad9de06b443582e2346c5fddcfb9b3b55ab03e562a5c0c4c158a,28,1.4928571428571427)
(fd9709ae2b08802a0cf32aa1971dd29c0de7c8b4be3cc07a1cb968fe2405ed5,28,1.4785714285714284)
(882faaed944cce28b59a882b46075b76dad00b0580d3012800a19a25ae9b3221,28,1.4000000000000001)
(3a2e34b3ebca5885323eb7b26d74eab2436cab04d22897c099635d550c9e9201,22,1.3090909090909089)
(df8239ca5372ac918ccc321ea281f93ca096766e88f3783471b81c7edc720b63,28,1.3071428571428572)
(41a9ffb21135b6c0a3007931bdf1ec944ab4434f1435924a008aa9e3a3ec15f,26,1.2769230769230768)
(9bb1f24866e6e0ab55847f85f87829cd4f7be72d1d341c9e514f14602e4e30d4,28,1.2642857142857142)
(075f4288380a972f084731c23f3ae382165107e4c5a2a2cd85363d3a96046fed,24,1.1333333333333333)
(e03396384ee196698bfc8bc0e21b2af1f2d72950c2d505a6be57261f1a2b6634,28,1.0357142857142858)
(84c39a5f9e43f9068fb7d4de358f0bbf67c90947463da8fe264c40c6895502e6,22,0.8909090909090901)
(5e66c5f5200f0426105d3639378ede436e1b0611b183df366fd42b5b4b3e7bac,28,0.5785714285714285)
(9582fbc05cf3288085a5d745452cebc2255674776c23c0495cd6d6e852418a02,28,0.5285714285714286)
(7fd5a24c7f7cae6655cc5747682409abc28f5872fdd868bb033310ab07b1fa0c,24,0.0)
(14c1e6fdb35fc08eb0fe6b496924fdb2280c15bb2ab9279e4ca9d4c8d73a4e2,20,0.0)
(746fbb665bfd41bf0470020cf596bea17d648b383f88f991408c55c191059b59,8,0.0)
(56a2e07bec0c4250925b2bb8579ac06a309404e9d03d911627b986a2f8ad57a7,8,0.0)
2025-04-14 23:43:35,602 [main] INFO org.apache.pig.Main - Pig script completed in 5 seconds and 5 milliseconds (5005 ms)

real    0m6.705s
user    0m17.949s
sys      0m1.493s
python3 hadoop2/hadoop2/hadoop2-#B-Python3-tester-44-cd1muw /Docker/comp-2/046
```

Elapsed time:5.005 seconds

Query 2: Faculty-wise Summary of Attendance and Course Credit

Objective:

Calculate the number of students, average attendance percentage, and maximum course credit for each faculty by filtering for passed students.

Steps and Logic:

1. Joining Datasets:

- Join the **fact_table** and **grade_roster** on student and course identifiers.

2. Filtering:

- Filter the joined dataset for records where the exam result is **'Pass'** to focus on successful outcomes.

3. Grouping by Faculty:

- Group the filtered records by **faculty_name**.

4. Aggregation:

- Count distinct students per faculty.

- Compute the average attendance using the `average_attendance_percent` from the fact table.
- Determine the maximum course credit awarded for courses taught by each faculty.

5. Output:

- Dump the results and store them into `/output/Query-2`.

Pig Script Excerpt:

```
-- Join fact_table and grade_roster
joined_data = JOIN fact_table BY (member_id, course), grade_roster BY
(student_id, subject_code_name);

-- Filter for students who passed
filtered_data = FILTER joined_data BY grade_roster::exam_result == 'Pass';

-- Group by faculty_name and compute aggregates
grouped_data = GROUP filtered_data BY grade_roster::faculty_name;
result = FOREACH grouped_data {
    unique_students = DISTINCT filtered_data.grade_roster::student_id;
    GENERATE group AS faculty_name,
              COUNT(unique_students) AS num_students,
              AVG(filtered_data.fact_table::average_attendance_percent) AS
avg_attendance,
              MAX(filtered_data.grade_roster::course_credit) AS
max_course_credit;
}

DUMP result;
STORE result INTO '/output/Query-2' USING PigStorage(',');
```

Image

```

2025-04-14 23:53:13.086 [main] WARN org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Encountered Warning ACCESSING_NON_EXISTENT_FIELD 30 tti
2025-04-14 23:53:13.086 [main] WARN org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Encountered Warning FIELD_DISCARDED_TYPE_CONVERSION_FAIL
s).
2025-04-14 23:53:13.086 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
2025-04-14 23:53:13.089 [main] WARN org.apache.pig.data.SchemaTupleBackend - SchemaTupleBackend has already been initialized
2025-04-14 23:53:13.091 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input files to process : 1
2025-04-14 23:53:13.091 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
(Vinu E V,59,87.05762688588288,4)
(V Sridhar,313,83.28612771888689,4)
(S Malapaka,166,80.00903669610081,4)
(Sachit Rao,150,74.71743122415805,4)
(Tulika Saha,120,73.93666648864746,2)
(Uttam Kumar,2,28.0,4)
(Vinod Reddy,5,67.04000091552734,4)
(Viswanath G,145,85.38620755425815,4)
(Nanditha Rao,42,66.87857182820638,4)
(Priyanka Das,6,77.18333371480306,4)
(Sakshi Arora,30,73.76666717529297,4)
(Thangaraju B,149,92.32364866218052,4)
(Jyotsna Bapat,2,97.20000076293945,4)
(Preeti Mudliar,33,80.20000053174567,4)
(Srinivas Vivek,198,77.55353602014407,4)
(Sushree Behera,4,81.82500076293945,4)
(Kurian Polachan,91,86.91978018624442,4)
(Priyanka Sharma,280,66.44857121195112,2)
(Ashish Choudhury,6,80.73333358764648,4)
(Manisha Kulkarni,119,76.56722676253118,4)
(Meenakshi D Souza,3,86.26666514078777,4)
(Srinath Srinivasa,3,88.9000015258789,4)
(Amit Chattopadhyay,159,84.39371070621898,4)
(Prof. Amrita Mishra,120,79.95333296457926,4)
(Badrinath Ramamurthy,120,87.22250073750814,2)
(G Srinivasa Raghavan,4,88.67499923706055,4)
(Jaya Sreevalsan Nair,1,70.80000305175781,4)
(Pillalamarri Sridhar,160,80.70999963283539,4)
(Sujit Kumar Chakrabarti,160,86.43625028133393,2)
(Karthikeyan Vaidyanathan,1,85.69999694824219,4)
2025-04-14 23:53:13.136 [main] INFO org.apache.pig.Main - Pig script completed in 4 seconds and 445 milliseconds (4445 ms)

real    0m6.143s
user    0m15.457s
sys     0m1.430s

```

Elapsed time:4.445 seconds

Query 3: Identify Low Attendance (Below 75%) per Student-Course

Objective:

Determine the attendance percentage for each student in each course and identify those records where attendance is below 75%.

Steps and Logic:

1. Joining Datasets:

- Join **fact_table** with **grade_roster** on matching student and course identifiers.

2. Grouping Data:

- Group the joined data by both **student_id** and **subject_code_name** to work at the granularity of each student's course.

3. Attendance Calculation:

- Compute total classes attended and absent for each group.
- Calculate the overall attendance percentage using the formula:

$$(\text{attendance_percentage}) = \frac{\text{total attended}}{\text{total attended} + \text{total absent}} \times 100$$

4. Filtering:

- Filter out groups where the attendance percentage is less than 75%.

5. Output:

- Dump and store the final filtered output into **/output/Query-3**.

Pig Script Excerpt:

```
-- Join fact_table and grade_roster
joined_data = JOIN fact_table BY (member_id, course), grade_roster BY
(student_id, subject_code_name);

-- Group by student_id and subject_code_name
grouped_data = GROUP joined_data BY (grade_roster::student_id,
grade_roster::subject_code_name);

-- Calculate attendance metrics per group
attendance_data = FOREACH grouped_data {
    total_attended =
SUM(joined_data.fact_table::number_of_classes_attended);
    total_absent = SUM(joined_data.fact_table::number_of_classes_absent);
    attendance_percentage = (total_attended * 100.0) / (total_attended +
total_absent);
    GENERATE FLATTEN(group) AS (student_id, course),
        total_attended AS total_classes_attended,
        total_absent AS total_classes_absent,
        attendance_percentage AS overall_attendance_percentage;
}

-- Filter groups with attendance below 75%
filtered_attendance = FILTER attendance_data BY
overall_attendance_percentage < 75;

DUMP filtered_attendance;
STORE filtered_attendance INTO '/output/Query-3' USING PigStorage(',');
```

Image

```
(f9746d5926e1ef8be988f4a01b8897189476a4792deee63c7aa37e2d31b862a3,EGC 112/Programming 1B (Python Programming),8,4,66.66666666666667)
(fa97ea0f7b79d3347a03f5cd5c5e96188d59f7e7098a0cec26b28d2f804fcf205,CSE 511/Algorithms,180,63,74.07407407407408)
(fbd0443bf1e0d231601b6aff94a29877222acac5946506425863c35151df2084,GNL 101/English,5,5,50.0)
(fc1e3958bf58979da2cd0fd53a5a62ba0377eb11aeb44e08b2ea5f37cc2f2fb,AIM 512/Mathematics for Machine Learning,1,1,50.0)
(fc43072bf0449e0f4f3743a9fb44d63507c0444bf6db7440443111fb0f406bce,GNL 101/English,3,8,27.272727272727273)
(fc43072bf0449e0f4f3743a9fb44d63507c0444bf6db7440443111fb0f406bce,AMS 103/Calculus,68,64,51.515151515151516)
(fc43072bf0449e0f4f3743a9fb44d63507c0444bf6db7440443111fb0f406bce,HSS 111/Economics-1,7,6,53.84615384615385)
(fc43072bf0449e0f4f3743a9fb44d63507c0444bf6db7440443111fb0f406bce,EGC 102/Digital Design,22,12,64.70588235294117)
(fc43072bf0449e0f4f3743a9fb44d63507c0444bf6db7440443111fb0f406bce,AMS 101/Probability & Statistics,68,32,68.0)
(fc43072bf0449e0f4f3743a9fb44d63507c0444bf6db7440443111fb0f406bce,EGC 111/Programming 1A (C Programming),9,7,56.25)
(fc43072bf0449e0f4f3743a9fb44d63507c0444bf6db7440443111fb0f406bce,EGC 112/Programming 1B (Python Programming),8,3,72.72727272727273)
(fc4535a76a081757ff741a0cf4f9aef52866e36e06aacc43239945bd0cca113c,GNL 101/English,4,7,36.36363636363637)
(fc4535a76a081757ff741a0cf4f9aef52866e36e06aacc43239945bd0cca113c,AMS 103/Calculus,96,36,72.72727272727273)
(fc4535a76a081757ff741a0cf4f9aef52866e36e06aacc43239945bd0cca113c,AMS 101/Probability & Statistics,72,28,72.0)
(fc4535a76a081757ff741a0cf4f9aef52866e36e06aacc43239945bd0cca113c,EGC 111/Programming 1A (C Programming),11,5,68.75)
(fc4535a76a081757ff741a0cf4f9aef52866e36e06aacc43239945bd0cca113c,EGC 112/Programming 1B (Python Programming),8,3,72.72727272727273)
(fc5f93239ec1b27fddbf7174a1f68e953d57e0b86e3c910135d02658a01a26ed,GNL 101/English,7,3,70.0)
(fcfa55660b5d441de2ef2e9b0b5b18c33a3f4853acdd231fea1eddd58dcd1ee,GNL 101/English,7,3,70.0)
(fcfbf656fb89ac105f2d0a8393c61f314a8449184a2f72349eef90b477c6c37b,VLS 502/Analog CMOS VLSI Design,13,9,59.09090909090909)
(fcfbf656fb89ac105f2d0a8393c61f314a8449184a2f72349eef90b477c6c37b,VLS 505/System design with FPGA,4,5,44.44444444444444)
(fcfbf656fb89ac105f2d0a8393c61f314a8449184a2f72349eef90b477c6c37b,VLS 864/Embedded Systems Design,16,8,66.66666666666667)
(fd9709ae2b0802a0cf32aa1971dd29c0de7c8b4be3cc07a1cb968fe2405ed5,EGC 112/Programming 1B (Python Programming),7,4,63.63636363636363)
(fdb1bf0b3ff8d8048103388f108794de4164bbe8dbdf7d898a6036965cc2f292,GNL 101/English,7,4,63.63636363636363)
(fdb1bf0b3ff8d8048103388f108794de4164bbe8dbdf7d898a6036965cc2f292,AMS 103/Calculus,92,40,69.6969696969697)
(fdb1bf0b3ff8d8048103388f108794de4164bbe8dbdf7d898a6036965cc2f292,EGC 102/Digital Design,25,9,73.52941176470588)
(fdb1bf0b3ff8d8048103388f108794de4164bbe8dbdf7d898a6036965cc2f292,AMS 101/Probability & Statistics,68,32,68.0)
(fdb1bf0b3ff8d8048103388f108794de4164bbe8dbdf7d898a6036965cc2f292,EGC 112/Programming 1B (Python Programming),7,4,63.63636363636363)
(fe6cadccebbf5892a3583e6ec13530f2e6ea7c6c75a90fccc9a2645e7200033,GNL 101/English,0,1,0.0)
(fe6cadccebbf5892a3583e6ec13530f2e6ea7c6c75a90fccc9a2645e7200033,AMS 103/Calculus,48,56,46.15384615384615)
(fe6cadccebbf5892a3583e6ec13530f2e6ea7c6c75a90fccc9a2645e7200033,HSS 111/Economics-1,2,10,16.666666666666668)
(fe6cadccebbf5892a3583e6ec13530f2e6ea7c6c75a90fccc9a2645e7200033,EGC 102/Digital Design,17,10,62.96296296296296)
(fe6cadccebbf5892a3583e6ec13530f2e6ea7c6c75a90fccc9a2645e7200033,AMS 101/Probability & Statistics,40,28,58.8235294117647)
(fedafcd150b9a17932760554a0ec9208266957a49da49214f4f9c7e1776f340d,GNL 101/English,7,3,70.0)
(ff6350e8fa8dc6311d81990d463738796ae3eb5cb545a29edad662c492864cbbf,VLS 505/System design with FPGA,6,3,66.66666666666667)
(ffba274d8a68b64e06980a5d807a0057faa389d2c7a5857424d47d9608e8c434,AIM 511/Machine Learning,0,4,0.0)
(ff3d0802fbf6b6c4020303b73c54bcef8c8e9c4b5db7108ac2c8f9b206f0f177,GNL 101/English,7,3,70.0)
(ff3d0802fbf6b6c4020303b73c54bcef8c8e9c4b5db7108ac2c8f9b206f0f177,EGC 111/Programming 1A (C Programming),32,28,53.33333333333333)
2025-04-14 23:54:24,976 [main] INFO org.apache.pig.Main - Pig script completed in 4 seconds and 536 milliseconds (4536 ms)

real    0m6.261s
user    0m16.118s
sys     0m1.387s
chachau: chandak@chachau: chandak #B: Build 1149, 2025-04-14 23:54:24, 976 [main] INFO org.apache.pig.Main - Pig script completed in 4 seconds and 536 milliseconds (4536 ms)
```

Elapsed time:4.536 seconds

Comparison of Hive vs. Pig

Table 1: Comparison of HiveQL Query Execution and Pig Query Execution

Query	Hive Without Bucketing(seconds)	Hive With Bucketing (seconds)	Pig performance (seconds)
Query 1	6.54	1.643	5.005
Query 2	0.912	0.793	4.445
Query 3	1.23	1.018	4.536

1. Installation & Setup

- **Hive:**
 - Typically involves setting up a Hive metastore along with Hadoop.
 - More components (HiveServer2, Metastore, etc.) need to be configured.
 - Can be complex to install and manage, especially in a production environment.
- **Pig:**
 - Generally easier to install and lightweight.
 - Runs as a single script without the need for a separate metastore.
 - Quick to set up on local mode or within a Hadoop cluster.

2. Query Language & Ease of Writing

- **Hive:**
 - Uses a SQL-like language (HiveQL) that is familiar to users with a relational database background.
 - Declarative queries make it easier for those accustomed to SQL.
 - Built-in functions and windowing can make complex queries simpler.
- **Pig:**
 - Uses a scripting language called Pig Latin, which is procedural.
 - Offers more flexibility and control when writing data transformation logic.
 - Can be easier for iterative data processing tasks, but may require more lines of code for similar SQL operations.

3. Query Performance & Optimization

- **Hive:**
 - Optimized for complex, long-running queries over large datasets.
 - Supports indexing, partitioning, and bucketing, which can significantly improve query performance when properly tuned.
 - More suitable for batch processing analytical queries.
- **Pig:**
 - Also handles large datasets but can be more efficient for ETL tasks and transformations.
 - Performance can be comparable to Hive for many transformation operations; however, highly optimized **Hive queries may outperform Pig on complex aggregations.**
 - Less emphasis on indexing and more on user-defined optimizations via scripting logic.

4. Suitability & Use Cases

- **Hive:**

- Best suited for analysts comfortable with SQL.
- Ideal for ad hoc queries and reporting where the data schema is well-defined.
- Strong integration with BI tools and reporting systems.
- **Pig:**
 - Excellent for ETL workflows and data processing pipelines.
 - Preferred when you need fine-grained control over data transformations.
 - Often used in scenarios where rapid prototyping of data flows is required.

5. Community & Ecosystem

- **Hive:**
 - Widely adopted in enterprises, with robust community support and integration with many Hadoop components.
 - Part of the broader SQL-on-Hadoop ecosystem.
- **Pig:**
 - Once very popular for data processing tasks, but usage has decreased in favor of Spark and other processing frameworks.
 - Still a viable option for specific transformation-heavy workflows.

In the table as well, we can see that Hive outperforms Pig on almost all queries in terms of time, and when Hive is optimised with partitioning and bucketing, it always outperforms Pig.

Also, it was a little difficult to write pig scripts since it is easier to write SQL statements, which is similar to Hive queries.