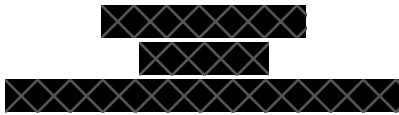# COMP30520 Practical 6

November 2020

# Exercise 1

## Exercise 1 - Question 1

```
1  movie_data = LOAD 'hdfs://localhost:8020/user/Movie_Log'
2      USING PigStorage(',')
3      AS (id: int,
4      rate_movie: chararray, rating: int,
5      completed_movie: chararray, genre: chararray,
6      pause_movie: chararray, start_movie: chararray,
7      browse_movie: chararray, list_movie: chararray,
8      search_movie: chararray, login: chararray,
9      logout: chararray, incomplete_movie: chararray,
10     purchase_movie: chararray);
11
12 DUMP movie_data;
```

## Exercise 1 - Question 2

```
1  rate_group = GROUP movie_data BY rating;
2  DUMP rate_group
```

## Exercise 1 - Question 3

```
1  DESCRIBE rate_group;
```

The DESCRIBE function shows use the Schema used for the Data we Describe. This Schema shows what the group is made of and the datatypes of the variables inside the Schema.

## Exercise 1 - Question 4

```
1  rate_group = FOREACH rate_group GENERATE group AS rating, movie_data.rate_movie;
```

## Exercise 1 - Question 5

```
1  grpd = GROUP movie_data BY genre;
2  max_rtd = FOREACH grpd {
3          ordrd = ORDER movie_data BY rating DESC;
4          top = LIMIT ordrd 1;
5          GENERATE top;
6  }
7  DUMP max_rtd;
```

## Exercise 1 - Question 6

```
1  selected_clicks = FOREACH movie_data GENERATE start_movie, browse_movie, completed_movie,
       purchased_movie;
```

# Exercise 2

## Exercise 2 - Question 1

```
1  student_data = LOAD 'hdfs://localhost:8020/user/students.csv'
2      USING PigStorage(',')
3      AS (id: int, first_name: chararray,
4          last_name: chararray, student_id: int,
5          email: chararray, stage: int,
6          degree: chararray);
7
8  DESCRIBE student_data
```

## Exercise 2 - Question 2

```
1  students_details = LOAD 'hdfs://localhost:8020/user/students.csv' USING PigStorage(',');
2  students_details = FOREACH students_details generate $3 as student_id,
3                                        $1 as first_name,
4                                        $2 as last_name,
5                                        $4 as email,
6                                        $5 as stage,
7                                        $6 as degree;
8
9  STORE renamed INTO 'hdfs://localhost:8020/pig_Output' USING PigStorage(',');
```

```
shell> hadoop fs -getmerge /user/pigoutput ./students_details.csv
```

## Exercise 2 - Question 3

```
1  student_data = LOAD 'hdfs://localhost:8020/user/students_attendance.csv'
2      USING PigStorage(',')
3      AS (id: int, student_id: int, hours_attended: int);
4
5  DESCRIBE student_data;
6
7  /*find all students that have attended*/
8  SA_details = FILTER student_data BY hours_attended > 1;
9  STORE SA_details INTO 'hdfs://localhost:8020/pig_OutputSA' USING PigStorage(',');
```

```
shell> hadoop fs -getmerge /user/pigoutput ./students_attendance.csv
```

## Exercise 2 - Question 4

```
1  attendance = LOAD 'hdfs://localhost:8020/user/students_attendance.csv' USING PigStorage(',')  AS (
       id: int, student_id: int, hours_attended: int);
2  data = GROUP attendance BY student_id;
3  result = FOREACH data GENERATE group, SUM(attendance.hours_attended);
4  DUMP result;
```

## Exercise 2 - Question 5

```
1  student_data = LOAD 'hdfs://localhost:8020/user/students.csv'
2      USING PigStorage(',')
3      AS (id: int, first_name: chararray,
4          last_name: chararray, student_id: int,
5          email: chararray, stage: int,
6          degree: chararray);
7
8  attendance = LOAD 'hdfs://localhost:8020/user/students_attendance.csv'
9      USING PigStorage(',')
10     AS (id: int, student_id: int, hours_attended: int);
11
12 joined = JOIN student_data BY student_id LEFT OUTER, attendance BY student_id;
13
14 result = FOREACH joined GENERATE student_data::student_id AS student_id,
15     student_id::first_name AS name,
16     attendance::hours_attended AS hours_attended;
```

## Exercise 2 - Question 6

```
1  DUMP result;
```

And

```
1  ILLUSTRATE result;
```