# Report
## Assignment 2 - MySQL

**Group**: 10
**Students**: Anna Sofie Lunde, Astrid Kristine Ragnhildsdatter Bakken and Maria Lundin Brenna

## Introduction

The task was to create tables, clean and insert data into the tables, and then query the database using Python and MySQL. The data we worked with was from Microsoft's Geolife GPS Trajectory dataset. The dataset contains data on users' outdoor movements and our task was inspired by the workout application Strava.

We mostly worked together during the creation of the tables and the data cleaning as we considered this part as crucial for the project. After the data was inserted into the database, each group member got responsibility for some queries. When each query had a draft, we discussed the different drafts and agreed on the final version of the query.

Link to our Git repository: https://github.com/annalunde/TDT4225-Assignment-2

# Results

## *Part 1*
User – top 10 rows

| id  | has_labels |
|-----|------------|
| 000 | 0          |
| 001 | 0          |
| 002 | 0          |
| 003 | 0          |
| 004 | 0          |
| 005 | 0          |
| 006 | 0          |
| 007 | 0          |
| 008 | 0          |
| 009 | 0          |

Activity – top 10 rows

| id | user_id | transportation_mode | start_date_time     | end_date_time       |
|----|---------|---------------------|---------------------|---------------------|
| 1  | 135     | None                | 2009-01-03 01:21:34 | 2009-01-03 05:40:31 |
| 2  | 135     | None                | 2009-01-02 04:31:27 | 2009-01-02 04:41:05 |
| 3  | 135     | None                | 2009-01-27 03:00:04 | 2009-01-27 04:50:32 |
| 4  | 135     | None                | 2009-01-10 01:19:47 | 2009-01-10 04:42:47 |
| 5  | 135     | None                | 2009-01-14 12:17:57 | 2009-01-14 12:30:53 |
| 6  | 135     | None                | 2009-01-12 01:41:22 | 2009-01-12 02:14:01 |
| 7  | 135     | None                | 2008-12-24 14:42:07 | 2008-12-24 15:26:45 |
| 8  | 135     | None                | 2008-12-28 10:36:05 | 2008-12-28 12:19:32 |
| 9  | 132     | None                | 2010-02-15 10:56:35 | 2010-02-15 12:22:33 |
| 10 | 132     | None                | 2010-04-30 23:38:01 | 2010-05-01 00:35:31 |

TrackPoint – top 10 rows

| id | activity_id | lat     | lon   | altitude | date_days | date_time           |
|----|-------------|---------|-------|----------|-----------|---------------------|
| 1  | 1           | 39.9743 | 116.4 | 492      | 39816.1   | 2009-01-03 01:21:34 |
| 2  | 1           | 39.9743 | 116.4 | 492      | 39816.1   | 2009-01-03 01:21:35 |
| 3  | 1           | 39.9743 | 116.4 | 492      | 39816.1   | 2009-01-03 01:21:36 |
| 4  | 1           | 39.9743 | 116.4 | 492      | 39816.1   | 2009-01-03 01:21:38 |
| 5  | 1           | 39.9744 | 116.4 | 491      | 39816.1   | 2009-01-03 01:21:39 |
| 6  | 1           | 39.9744 | 116.4 | 491      | 39816.1   | 2009-01-03 01:21:42 |
| 7  | 1           | 39.9744 | 116.4 | 491      | 39816.1   | 2009-01-03 01:21:46 |
| 8  | 1           | 39.9745 | 116.4 | 491      | 39816.1   | 2009-01-03 01:21:51 |
| 9  | 1           | 39.9745 | 116.4 | 490      | 39816.1   | 2009-01-03 01:21:56 |
| 10 | 1           | 39.9745 | 116.4 | 489      | 39816.1   | 2009-01-03 01:22:01 |

*Part 2*

Note: The assumptions for the different queries are listed under *Discussion*

*Query 1:*

Number of users: 182
Number of activities: 16048
Number of trackpoints: 9681756

| NumUsers | NumActivities | NumTrackpoints |
|---|---|---|
| 182 | 16048 | 9681756 |

*Query 2:*

Maximum number of activities per user: 2102
Minimum number of activities per user: 1
Average number of activities per user: 92.763

| Maximum | Minimum | Average |
|---|---|---|
| 2102 | 1 | 92.763 |

*Query 3:*

Top 10 users with the highest number of activities:

| user_id | Count |
|---|---|
| 128 | 2102 |
| 153 | 1793 |
| 025 | 715 |
| 163 | 704 |
| 062 | 691 |
| 144 | 563 |
| 041 | 399 |
| 085 | 364 |
| 004 | 346 |
| 140 | 345 |

## Query 4:

The number of activities for users, where the activity is started in one day and ended in the next day:

| user_id | NumActivities | user_id | NumActivities | user_id | NumActivities |
|---------|---------------|---------|---------------|---------|---------------|
| 000 | 5 | 041 | 14 | 111 | 4 |
| 001 | 19 | 042 | 1 | 112 | 2 |
| 002 | 12 | 043 | 2 | 113 | 2 |
| 003 | 6 | 044 | 13 | 115 | 28 |
| 004 | 28 | 050 | 1 | 118 | 1 |
| 005 | 1 | 051 | 2 | 124 | 1 |
| 006 | 1 | 052 | 5 | 125 | 3 |
| 007 | 1 | 053 | 1 | 126 | 3 |
| 010 | 19 | 057 | 6 | 128 | 33 |
| 011 | 4 | 058 | 1 | 129 | 1 |
| 013 | 36 | 061 | 6 | 132 | 1 |
| 014 | 30 | 062 | 4 | 134 | 2 |
| 015 | 18 | 067 | 2 | 138 | 1 |
| 016 | 1 | 068 | 59 | 139 | 2 |
| 017 | 16 | 069 | 1 | 140 | 5 |
| 018 | 3 | 070 | 2 | 142 | 2 |
| 019 | 1 | 071 | 2 | 144 | 48 |
| 020 | 1 | 074 | 1 | 146 | 2 |
| 021 | 2 | 076 | 1 | 147 | 2 |
| 022 | 7 | 081 | 1 | 150 | 6 |
| 024 | 5 | 082 | 2 | 153 | 236 |
| 025 | 9 | 083 | 3 | 155 | 2 |
| 026 | 1 | 084 | 35 | 157 | 1 |
| 027 | 1 | 085 | 32 | 163 | 24 |
| 028 | 19 | 088 | 2 | 167 | 6 |
| 029 | 18 | 091 | 2 | 168 | 8 |
| 030 | 12 | 092 | 4 | 172 | 1 |
| 032 | 6 | 094 | 6 | 174 | 2 |
| 035 | 11 | 095 | 1 | 175 | 1 |
| 036 | 3 | 099 | 1 | | |
| 037 | 15 | 100 | 1 | | |
| 038 | 7 | 101 | 1 | | |
| 039 | 37 | 103 | 3 | | |
| | | 104 | 5 | | |
| | | 106 | 2 | | |
| | | 108 | 1 | | |

### Query 5:

Activities that are registered multiple times: Zero results

```
user_id    transportation_mode    start_date_time    end_date_time    NumDuplicates
---------  ---------------------  -----------------  ---------------  ----------------


-----------------------------------------------
```

### Query 6:

Have tried to run query 6, but the query is too heavy, so no results are reported.

```python
def query_six(self, table_name_activities, table_name_trackpoints):
    """
    Find the number of users which have been close to each other in time and
    space (Covid-19 tracking). Close is defined as the same minute (60 seconds)
    and space (100 meters).
    """

    query = "SELECT t1.user_id, t1.lat, t1.lon, t2.user_id, t2.lat, t2.lon " \
            "FROM (SELECT user_id, lat, lon, date_time FROM %s inner join %s on Activity.id=TrackPoint.activity_id) as t1, " \
            "(SELECT user_id, lat, lon, date_time FROM Activity inner join TrackPoint on Activity.id=TrackPoint.activity_id) as t2 " \
            "where t1.user_id != t2.user_id " \
            "AND ABS(TIMESTAMPDIFF(SECOND,t1.date_time, t2.date_time)) <= 60" \

    self.cursor.execute(query % (table_name_activities,
                                 table_name_trackpoints))
    rows = self.cursor.fetchall()
    print(tabulate(rows, headers=self.cursor.column_names))

    user_dict = dict()
    for row in rows:
        if (haversine((row[1], row[2]), (row[4], row[5]), unit="km") <= 0.1):
            if row[0] in user_dict:
                user_dict[row[0]].append(row[3])
            else:
                user_dict[row[0]] = [row[3]]
    users = 0
    for value in users_dict.values():
        users += len(value)
    users = users/2
    print(users)
    return users
```

### Query 7:

All users that have never taken a taxi:

| user_id |
|---------|
| 010 |
| 020 |
| 021 |
| 052 |
| 056 |
| 058 |
| 060 |
| 062 |
| 064 |
| 065 |
| 067 |
| 069 |
| 073 |
| 075 |
| 076 |
| 078 |
| 080 |
| 081 |
| 082 |
| 084 |
| 085 |
| 086 |
| 087 |
| 089 |
| 091 |
| 092 |
| 097 |
| 101 |
| 102 |
| 107 |
| 108 |
| 112 |
| 115 |
| 117 |
| 125 |
| 126 |
| 128 |
| 136 |
| 138 |
| 139 |
| 144 |
| 153 |
| 161 |
| 163 |
| 167 |
| 175 |

*Query 8:*

All types of transportation modes and how many distinct users have used the different transportation modes:

```
TransportationMode       NumDistinctUsers
--------------------     ------------------
airplane                                1
bike                                   19
boat                                    1
bus                                    13
car                                     8
run                                     1
subway                                  4
taxi                                   10
train                                   2
walk                                   34
```

*Query 9:*

a) Year and month with the most activities: November 2008

```
Year     Month     ActivityCount
------   -------   ---------------
2008      11              1006
```

b) The user that had the most activities in November 2008, and the number of recorded hours: User 062 with 130 activities and 7 hours recorded

The user with the most activities in November 2008 (User 062 – 7 hours) does not have more hours recorded than the user with the second most activities (User 128 – 34 hours).

```
user_id     ActivityCount     HoursActive
--------   ----------------   -------------
    062                130               7
    128                 75              34
```

*Query 10:*

Total distance (in km) walked in 2008 by user with id=112:
115.47465961507991 km

*Query 11:*

Find the top 20 users who have gained the most altitude:

```
user_id    MetersGained
--------   ---------------
    128         650887
    153         554969
    004         332036
    041         240758
    003         233664
    085         217642
    163         205264
    062         181692
    144         179457
    030         175680
    039         146704
    084         131161
    000         121505
    002         115063
    167         112973
    025         109148
    037          99220.9
    140          94838.8
    126          83024.2
    017          62566.3
```

*Query 12:*

All users who have invalid activities, and the number of invalid activities per user:

| user_id | NumInvalid | user_id | NumInvalid | user_id | NumInvalid |
|---------|-----------|---------|-----------|---------|-----------|
| 000 | 101 | 031 | 3 | 061 | 12 |
| 001 | 45 | 032 | 12 | 062 | 249 |
| 002 | 98 | 033 | 2 | 063 | 8 |
| 003 | 179 | 034 | 88 | 064 | 7 |
| 004 | 219 | 035 | 23 | 065 | 26 |
| 005 | 45 | 036 | 34 | 066 | 6 |
| 006 | 17 | 037 | 100 | 067 | 33 |
| 007 | 30 | 038 | 58 | 068 | 139 |
| 008 | 16 | 039 | 147 | 069 | 6 |
| 009 | 31 | 040 | 17 | 070 | 5 |
| 010 | 50 | 041 | 201 | 071 | 29 |
| 011 | 32 | 042 | 55 | 072 | 2 |
| 012 | 43 | 043 | 21 | 073 | 18 |
| 013 | 29 | 044 | 32 | 074 | 19 |
| 014 | 118 | 045 | 7 | 075 | 6 |
| 015 | 46 | 046 | 13 | 076 | 8 |
| 016 | 20 | 047 | 6 | 077 | 3 |
| 017 | 129 | 048 | 1 | 078 | 19 |
| 018 | 27 | 050 | 8 | 079 | 2 |
| 019 | 31 | 051 | 36 | 080 | 6 |
| 020 | 20 | 052 | 44 | 081 | 16 |
| 021 | 7 | 053 | 7 | 082 | 27 |
| 022 | 55 | 054 | 2 | 083 | 15 |
| 023 | 11 | 055 | 15 | 084 | 99 |
| 024 | 27 | 056 | 7 | 085 | 184 |
| 025 | 263 | 057 | 16 | 086 | 5 |
| 026 | 18 | 058 | 13 | 087 | 3 |
| 027 | 2 | 059 | 5 | 088 | 11 |
| 028 | 36 | 060 | 1 | 089 | 40 |
| 029 | 25 | | | 090 | 3 |
| 030 | 112 | | | | |

| user_id | NumInvalid | user_id | NumInvalid | user_id | NumInvalid |
|---------|-----------|---------|-----------|---------|-----------|
| 091 | 63 | 123 | 3 | 158 | 9 |
| 092 | 101 | 124 | 4 | 159 | 5 |
| 093 | 4 | 125 | 25 | 161 | 7 |
| 094 | 16 | 126 | 105 | 162 | 9 |
| 095 | 4 | 127 | 4 | 163 | 233 |
| 096 | 35 | 128 | 720 | 164 | 6 |
| 097 | 14 | 129 | 6 | 165 | 2 |
| 098 | 5 | 130 | 8 | 166 | 2 |
| 099 | 11 | 131 | 10 | 167 | 134 |
| 100 | 3 | 132 | 3 | 168 | 19 |
| 101 | 46 | 133 | 4 | 169 | 9 |
| 102 | 13 | 134 | 31 | 170 | 2 |
| 103 | 24 | 135 | 5 | 171 | 3 |
| 104 | 97 | 136 | 6 | 172 | 9 |
| 105 | 9 | 138 | 10 | 173 | 5 |
| 106 | 3 | 139 | 12 | 174 | 54 |
| 107 | 1 | 140 | 86 | 175 | 4 |
| 108 | 5 | 141 | 1 | 176 | 8 |
| 109 | 3 | 142 | 52 | 179 | 28 |
| 110 | 17 | 144 | 157 | 180 | 2 |
| 111 | 26 | 145 | 5 | 181 | 14 |
| 112 | 67 | 146 | 7 | | |
| 113 | 1 | 147 | 30 | | |
| 114 | 3 | 150 | 16 | | |
| 115 | 58 | 151 | 1 | | |
| 117 | 3 | 152 | 2 | | |
| 118 | 3 | 153 | 557 | | |
| 119 | 22 | 154 | 14 | | |
| 121 | 4 | 155 | 30 | | |
| 122 | 6 | 157 | 9 | | |

## Discussion

*Part 1*
- To keep track of the activityID given in the database, the file name of the activity along with activityID were saved in a separate .txt-file. This .txt-file was used when data was inserted into the TrackPoint table, so that a trackpoint would be linked to the correct activity with the correct activityID.

*Part 2*
- In Query 4, we found the number of activities for users, where the activity is started in one day and ended in the next day (as mentioned on Piazza).
- Query 6 was very heavy, which led to it not returning any results.
- For Query 7, we assumed that we should only consider labelled activities, but that not all activities were to be labelled for that user to be considered as never have taken a taxi.
- For Query 9a, we assumed that the activities belonged to the year and month according to their start_date_time. It was deemed reasonable due to that it would be very complicated to keep track of the "border activities" - that started in one year/month and ended in another year/month. We also assumed that this would apply to few activities.
- A tip that was given in the assignment sheet was that variables in SQL might come in handy. For Query 9b, we decided to use the output for month and year from Query 9a instead, since this was an easier solution that we knew would provide the correct answer. However, it might have been a more elegant solution to use variables.

*Learning points*
- During this exercise we refreshed our knowledge within SQL and got practice in setting up and populating a database.
- We also learned the importance of knowing the structure of your dataset. We had to really get familiar with the data before being able to create tables, populate them, and write queries to extract the desired data.

## Feedback

The exercise was interesting and fun to execute. However, Query 6 was time consuming, and it would have been nice with more tips on how to test queries that are too heavy to be executed within a reasonable time frame.