

Assignment 4

Hadoop

Toon Van Craenendonck
toon.vancraenendonck@cs.kuleuven.be

Jessa Bekker
jessa.bekker@cs.kuleuven.be

Collaboration policy:

Projects are independent: no working together! You must come up with how to solve the problem independently. Do not discuss specifics of how you structure your solution, etc. You cannot share solution ideas, pseudocode, code, reports, etc. If you are unsure about the policy, ask the professor in charge or the TAs.

Before you start... note that you will have an exercise session on Hadoop on the 25th of April. If you would like to start working on this assignment before, check the Exercise 6 sheet for more details about working with Hadoop.

1 Datasets

For this assignment, you will be working on GPS tracks of taxis in San Francisco¹. We provide several datasets, which are briefly described in this section. More details about their format are given later, in the corresponding sections. The datasets are available online² and will be available on the clusters in the `/data` directory.

- **all.segments:** The complete GPS tracks (including intermediary points) from May 2008 to January 2011. Each record is a segment with two end points (start position, end position), two time stamps (start date, end date) and two taxi states (empty/full) for the beginning and for the end of the segment. Consecutive segments can be concatenated to form complete trips. The dataset contains about 306 million segments and weighs 28GiB.
- **2010.03.segments:** Follows the same specification as the previous dataset, but only contains the segments started or finished in March 2010. The dataset still contains about 19 million segments and weighs 1.2GiB.
- **2010.03.trips:** Contains trips constructed from the segments in the previous dataset. Each trip is simply represented by its two end points (i.e. there is no intermediary position).
- **Taxi_706.segments:** This dataset contains 19237 segments from taxi number 706, exclusively. This dataset is useful for debugging.

¹More details at cabspotting.org

²http://people.cs.kuleuven.be/~toon.vancraenendonck/bdap_files/

The datasets are provided (almost) as such: they are not sorted, and they contain errors and misformatted records. Dealing with this type of data is part of the assignment. All datasets will be available on the distributed file system (*DFS*). We recommend that you leave the largest dataset on the file system as there is no need to process it locally.

2 Trip length distribution

First, we are interested in computing a simple statistic: the distribution of trip lengths. We will compute this distribution for the trips in the `2010.03.trips` preprocessed dataset.

For this exercise only, we will assume that the trip distance is the distance between the two end points of the trip. This information is easy to compute from `2010.03.trips`, as it contains descriptions of the trips without the intermediary segments. In this dataset, each line has the following format (represented here on two lines):

```
<taxi-id> <start date> <start pos (lat)> <start pos (long)> ...  
...<end date> <end pos (lat)> <end pos (long)>
```

To compute the geographical distance between two coordinates, you can use a simple flat-surface formula³, which will give a reasonable approximation for this dataset because the distances are not too large (but remember that these formulæ are not always appropriate for larger distances⁴).

First, implement a simple algorithm in the language of your choice to compute the trip distance distribution. Then propose a MapReduce implementation that solves the same problem. Run it locally, then on the cluster, then compare your results. Which implementation is faster? Does it match your expectations? Include a description of your implementation and a discussion about efficiency in your report.

Plot the trip length distribution and interpret the results.

3 Computing airport ride revenue

A significant number of taxi rides pass through the San Francisco airport. Assume that taxi companies have to pay for an expensive license for this airport access. A company may then be interested in knowing exactly how much they earn from these airport rides, to know whether paying the license is actually worth it.

In this part of the assignment, you will be asked to give an estimate of the revenue coming from airport rides, based on the GPS tracking data. The estimate should be as accurate as possible, which means that you should consider all data, i.e. sampling is not an option.

This part of the assignment consists of two steps, which we describe next: (1) reconstructing trips from segments, (2) computing the revenue obtained from these trips.

3.1 Reconstructing trips

First, you will be required to reconstruct complete trips from the ride segments. The `.segments` files contain the complete GPS tracks decomposed into segments. A segment is simply a pair of geographical coordinates.

³http://en.wikipedia.org/wiki/Geographical_distance

⁴After all, the earth is not really flat.

The sampling rate is generally 1 minute, although there can be larger gaps. Each line from these datasets has the following format (represented here on two lines):

```
<taxi-id>, <start date>, <start pos (lat)>, <start pos (long)>, <start status> ...  
...<end date> <end pos (lat)> <end pos (long)> <end status>
```

Propose a design of a Map/Reduce application to construct trips. Depending on your implementation, you may or may not need multiple Map/Reduce jobs. Implement your application and test it locally on the `2010_03.segments` dataset.

Note on erroneous data-points: GPS points and recording devices are far from 100% reliable. Erroneous records will ultimately lead to erroneous results. One possible way to eliminate trips including erroneous data points, is to use a simple heuristic. For example, you can eliminate trips that include at least one segment with an average speed above 200km/h . You are free to come up with more elaborated ways of getting rid of erroneous trips as long as you discuss them in the report, possibly illustrating with examples of the erroneous trips. Plotting trips can be helpful to test the validity of your implementation in general. This can for example be done with the Google maps API (copying the link in the footnote in your browser should for example plot a simple fictitious trip between two points ⁵).

Once you have developed and tested your implementation on the sample, you should run it on the cluster on the complete dataset. However before you do so, you should think about the scalability of your approach. Execution times on Hadoop clusters are usually very variable and are thus not very relevant. Instead it is better to reason about efficiency in terms of number of input records and output records for the various components (mapper, combiner, reducer ...), and in terms of individual task complexity and maximum memory usage for each mapper and reducer. With the clusters you are given, how should you choose the number of mappers and the number of reducers? Based on this analysis, propose a set of changes to improve scalability of your approach. Test and observe the impact of your changes locally, and run your approach on the cluster to construct all trips.

3.2 Computing the revenue

In this step, you will use the output of the previous component to compute the total revenue obtained from airport trips. We consider airport trip rides as those that pass through a circle with the airport as center, and a radius of 1km. The airport is located at 37.62131° N, -122.37896° W. To calculate trip revenue, you can use a simple formula that combines a starting fee of \$3.5 with an additional \$1.71 per kilometer⁶.

Report the total revenue that has been earned from airport trips, and also make a plot that shows the evolution of this revenue over time.

4 If you want to go beyond ...

More interesting insights can be learned from this dataset, if you're interested in searching for more information feel free to do so and to comment in the report. As usual, any interesting insight about the problem or the data can potentially result in bonus points. (Nevertheless, you should focus on the previous sections first).

⁵ <https://maps.google.com/maps/api/staticmap?&size=640x640&markers=color:green%7Clabel:S%7C37.762573,-122.437477&markers=color:red%7Clabel:E%7C37.7452,-122.458076&path=geodesic:true|color:0x0000ff|weight:5|37.762573,-122.437477|37.7452,-122.458076>

⁶http://www.numbeo.com/taxi-fare/city_result.jsp?country=United+States&city=San+Francisco%2C+CA

5 (Very) Important remarks

- Here is a (non exhaustive) list of potential issues you may encounter when you're dealing with the datasets:
 - Date and timezone: to parse dates and compute trip durations, make sure you use the correct time zone (which can be fetched with `TimeZone.getTimeZone("America/San_Francisco")`).
 - Do not make wrong assumptions about the order or the format of the input record, your program should be robust enough to cope with broken records in the data without crashing.
- Do **NOT** run anything on the clusters before making sure it runs locally (at least on a smaller dataset). The clusters are shared resources.
- **Kill your jobs** on the cluster, if you think they are going to crash or run for too long. For that use `hadoop job -list` and `hadoop job -kill <job id>`. Not following this policy will have unimaginable consequences :)
- Scaling up to the large dataset is not trivial and may require some work. Make sure you have a decent implementation and a report for the smaller dataset (2010_03.segments) before working on the scalability on the large dataset. An implementation that works on the small dataset and a good report are worth more points than a complex undocumented implementation that runs on the larger dataset.

6 Report

Your report should be no more than four pages of text. Graphics, plots, tables, etc. do not count towards this total and you may include as many of them as you would like. Your report should address the following issues:

1. What approaches you used to identify incorrect data
2. A plot of the trip length distribution on the 2010_03.trips data, interpretation of the results, and a discussion of the run times for the different implementations
3. The total revenue that has been earned from airport trips, and also make a plot that shows the evolution of this revenue over time
4. A discussion about the numbers of mappers, reducers used for the revenue computation and their effect on scalability
5. A discussion of the mappers, reducers, combiners used for each task

7 Turning in your code and report

- The assignment should be handed in on Toledo before the **Friday the 27th of May**. As usual, there will be a 10% penalty per day, starting from the due day.
- You must upload an archive (.zip or .tar.gz) containing the following files
 - the report (as pdf)

- a runnable jar file called **Exercise1.jar** for the first exercise
- a runnable jar file called **Exercise2.jar** for the second exercise
- your source code in a sub directory called **src** (which may contain more subdirectories if necessary)
- a **README** including the command lines to compile and execute your code (on the cluster).

Good luck!