# Reproducibility review of: Process Analysis in Humanitarian Voluntary Geographic Information: the case of the HOT Tasking Manager

Rémy Decoupes [iD]

2024-05-27

## Reviewed paper

## Summary

The authors made their code available through a GitHub repository. The Readme.md file provides clear instructions for executing the entire code. It is worth noting that the input data is available provided that accounts are created with HOTOSM and Bunting Labs. The data retrieval process takes a long time. I initiated the retrieval process for 24 hours but only managed to receive one-third of the data. Therefore, the authors provided me with all the input data for my review.

I did not encounter any difficulties in executing the Python notebooks. As a result, I produced intermediate data (in 4 CSV files). Unfortunately, I encountered some blocking errors when executing the R Markdown (which is the notebook responsible for producing tables and figures of the manuscript).

Assisted by the authors, we discovered that the errors were caused by the content of the intermediate files I had generated. Strictly speaking, since the authors shared all their input data with me, I should not have encountered these errors.

To continue my review, the authors sent me their version of the intermediate files. Thanks to them, I was able to execute the entire R Markdown and generate all the tables and figures. At the time of finalizing

this report, I do not know why the generated intermediate files were different from those of the authors. I cannot affirm that the Python notebooks produce intermediate data of good quality.

# Reproducibility reviewer notes

The authors shared their code through a GitHub repository. The Readme.md file clearly explains the prerequisites and the steps to follow to reproduce the figures and tables. Four notebooks are provided, the 3 firsts in Python and the last one in R Markdown:

1. Download_data.ipynb
2. Build_density.ipynb
3. Data_preprocessing.ipynb
4. Process_Analysis.Rmd

What's more, to download the data, two accounts need to be created to activate API keys: OpenStreetMap and BuntingLags.

I ran the reproducibility review on Linux (Debian 10) with 32GB RAM.

## Python environment installation

First, I created a new Conda environment with dependencies as shared in the Readme.md file.

```
conda create -n agile-21 python=3.10.2 ipython pip ipykernel
conda activate agile-21

python ipykernel install --user --name=agile-21

pip install pandas==1.5.3
pip install geopandas==0.13.2
pip install numpy==1.25.2
pip install requests==2.31.0
pip install json==2.0.9
pip install ipywidgets==7.7.1
pip install tqdm==4.66.1
pip install utm==0.7.0
```

Then I ran the notebooks sequentially.

## 1. Download_data.ipynb

I followed the instructions to create the two tokens API from: - HOTOSM Tasking Manager API. I had to create first an OpenStreetMap account and authorize HOTOSM to access to my OpenStreetMap account. I faced an issue with Firefox, the redirection from OpenStreetMap to HOTOSM was not working but it worked with Chromium. - Buntinglabs: I just provided my e-mail address.

I create manually a directory `data` in which the notebook will save the collected data.

I had to stop my computer during the cell "Get project activities" after 20 hours of run. But when I re-start the notebook, the cell only downloaded the missing file (it was not trying to download already downloaded files). Unfortunately, even after 24 hours of execution, I had barely exceeded the 1/3 of the necessary data. So, to continue my review within the allotted time, I asked the authors to provide me with their input data.

The authors granted me access to all the datasets that this notebook retrieves, as I have an account and API key from HOT and Bunting Labs. This should be compliant with their terms of use.

## 2. Building Density

I had to setting a path to the output folder at the beginning of the notebook. At the end, it creates a csv file

## 3. Data preprocessing

- Missing installation of matplotlib: `pip install matplotlib`

Figure 1: agile-21-get-activity_too-long

- Error in the path to the file output_densities.csv in section 3. `Create file for regression analysis with task density ("regression.csv.csv)`, ll `densities=pd.read_csv("output_densities.csv")` should be `densities=pd.read_csv(data_folder + "output_densities.csv")` as data_folder containing output_densities.csv has been set in notebook `building_densitiy.ipynb`

At this end, this notebook created 4 csv files: - project.csv - initial_tasks.csv - contributors.csv - regression.csv

## 4. Process Analysis

As the notebook is a Rmarkdown, and in order to have the same R kernel as recommended by the authors, I used a docker image built by the Rocker Project. The image `rocker/rstudio:4.3.2` contains an RStudio server with R version 4.3.2. I mount my local directory (with all the code and data for reproducing this paper) into `/agile` inside the container.

```
sudo docker run -e ROOT=true -e PASSWORD=agile --rm -p 8787:8787 -v ./:/agile rocker/rstudio:4.3.2
```

When the Docker container is started, I connect to RStudio with a web browser: http://localhost:8787, with credentials as follows:

- login: rstudio
- password: agile

Then I installed the dependencies as listed in the Readme.md:

```
install.packages("bupaverse", version = "0.1.0")
install.packages("reshape2", version = "1.4.4")
install.packages("gt", version = "0.10.1")
install.packages("scales", version = "1.3.0")
install.packages("readr", version = "2.1.4")
install.packages("dplyr", version = "1.1.4")
install.packages("magrittr", version = "2.0.3")
install.packages("ggplot2", version = "3.4.4")
install.packages("Hmisc", version = "5.1-1")
install.packages("gamlss", version = "5.4-20")
```

Unfortunately I encountered several errors some of them were blocking:

- Missing installation of library:

```
Execution halted

No LaTeX installation detected (LaTeX is required to create PDF output). You should install a LaTeX distribution for your platform: https://www.latex-project.org/get/

  If you are not sure, you may install TinyTeX in R: tinytex::install_tinytex()

  Otherwise consider MiKTeX on Windows - http://miktex.org

  MacTeX on macOS - https://tug.org/mactex/
  (NOTE: Download with Safari rather than Chrome _strongly_ recommended)

  Linux: Use system package manager
```

So I installed tinytex:

```
tinytex::install_tinytex()
```

- In line 177, I got this error

```
Error in `group_by()`:
! Must group by variables found in `.data`.
 Column `mappingLevel` is not found.
Backtrace:
1. event_log_df %>% group_by(mappingLevel) %>% ...
4. dplyr:::group_by.data.frame(., mappingLevel)
```

It seems that in `event_log_df`, the correct name of the column is `mappingLevel.x` instead of `mappingLevel`.

I replaced by

```
mappingLevel <- event_log_df %>%  group_by(mappingLevel.x) %>% summarise(count = n_distinct(actionBy))
```

- In line 185, I got this error

```
Error in xtfrm.data.frame(x) : cannot xtfrm data frames
```

It seems that the dcast is missing a dimension (`data_pivot <- dcast(event_log_df, action ~ mappingLevel,value.var = "taskId", length)`). I could not find a fix but since `data_pivot` is not further used, I ignore this error

- In line 209, I got this error

```
Error in gamlss(percentage_area_covered_by_building ~ splits1 + invalidations1 +  :
  could not find function "gamlss"
```

I installed gamlss: `install.packages("gamlss")` and I had import it the cell in l.209: `library(gamlss)`.

- But then, I encountered this error that I could not bypass

```
Error in `contrasts<-`(`*tmp*`, value = contr.funs[1 + isOF[nn]]) :
  contrasts can be applied only to factors with 2 or more levels
```

As I was blocked by this issue, the Authors helped me to debug. They compared my intermediate CSV files generated by the three python notebooks. It appeared that they had significant differences affecting two file: regression.csv and initial_tasks that they were generated by the 3rd python notebook. In writing this report, I'm unsure where these differences might originate. Is it due to the input data, or perhaps there's an issue with the Python notebooks?

In order to test the generation of the figures and tables, I asked the authors to also send me the intermediate files. I was then able to execute the entire notebook without any errors.

Here in details the tables and figures I was able to reproduce

- Table 1:



Figure 2: Reproduction of table 1

- Table 3:



Figure 3: Reproduction of table 3

- Figure 7: Note that there are differences between those two graph, it's because I generated mine on a subset of data.
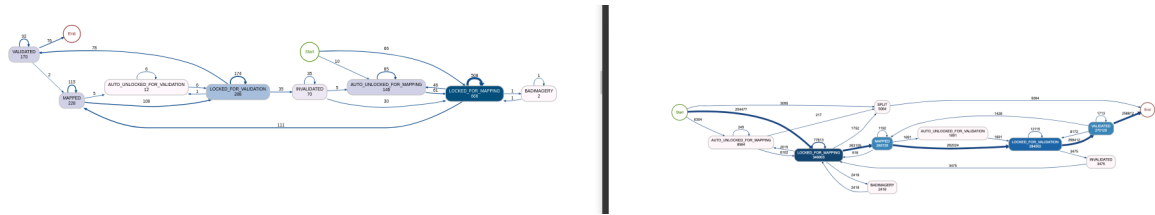


Figure 4: Reproduction of fig 7

- figure 8:
- Table 4: It seems that there is a little difference between the tables. This may be due to a difference in rounding or truncation.
- Table 5:
- Figure 9:
- Figure 10:

Figure 5: Reproduction of fig 8



Figure 6: Reproduction of Table 4



Figure 7: Reproduction of the 1rst part of Table 5



Figure 8: Reproduction of the 2nd part of Table 5

Figure 9. Time map of task states and transitions within the HOT-TM.

Figure 9: Reproduction of fig 9



Figure 10: Reproduction of fig 10