# The DEMIX Wine Contest Jupyter notebook

Carlos Henrique Grohmann

Universidade de São Paulo, Institute of Energy and Environment
Av. Prof. Luciano Gualberto, 1289, 05508-010
São Paulo, SP, Brazil
guano@usp.br

*Abstract*— **Jupyter Notebooks are the most widely used system for shareable and reproducible research. Within the DEMIX Working Group, a Python Jupyter Notebook was developed for the analysis and visualization of the DEMIX Wine Contest. It can be run on a user's local machine or in the cloud using the Google Colab platform. The notebook, as well as auxiliary files, are available in GitHub. This article briefly describes its characteristics and usage.**

## I.  INTRODUCTION

Jupyter Notebooks have become the most widely used system for shareable and reproducible research [1-3]. It supports a variety of programming languages, such as Julia, R, JavaScript, and C, allowing the user to create literate programming [4] documents combining code, text, and results with visualizations and other rich media [5].

Within the DEMIX Working Group [6], a Jupyter Notebook was developed for the analysis and visualization of the DEMIX Wine Contest. This article briefly describes its characteristics and usage. The reader is referred to [7] for the details of the DEMIX Wine Contest. In summary, the contest is designed to rank a set of Global DEMs based on a number of objective criteria (although subjective criteria could be used as well). Each criterion (such as RMSE) is used to rank the Global DEMs against a reference DEM. Then, a set of criteria is defined and the Friedman Test [8] is used to determine if the DEMs can or cannot be considered statistically different. If there are statistically significant differences among the set of DEMs, they are compared pairwise using the test by Dunn with Bonferroni correction [9] (a step called post-hoc analysis) and a final ranking is produced.

The notebook, as well as auxiliary files, are available in GitHub [10].

## II.  CHARACTERISTICS

The DEMIX Jupyter Notebook was developed based on Python version 3.10.x and the Pandas, qgrid, Numpy, Matplotlib, and Seaborn libraries [11-16].

As the idea behind the notebook was to provide a simple "interface" to analyze the DEMIX Wine Contest data, the statistical and plotting functions were implemented in the `demix_wine_functions.py` file, although the user can alter some plotting options, such as colors or symbols.

The notebook can be run on a user's local machine or in the cloud using the Google Colaboratory (Colab) platform. While running locally requires setting up a working python/jupyter environment, it allows for more flexibility in terms of files' location. Running it in Colab requires installing a specific version of qgrid, and downloading external files from GitHub and Zenodo. All the instructions necessary for the user to run the notebook are included as comments within code cells or as rich text.

## III.  USAGE

As input the Jupyter notebook takes the Wine Contest GIS database [17] produced by MICRODEM [18]. The GIS database contains signed (mean, median) and unsigned values (RMSE, LE90, MAE) of the differences of elevation, slope and roughness between the Global DEMs and reference DEMs, but only the unsigned values are read for the Wine Contest. The signed values are used to produce plots that help the user to explore and understand the set of data being analyzed.

Before exploring the database, the user can define the values for tolerances, which will impact the occurrence of ties between the DEMs in the rankings.

**Display Database with qgrid**

Any filtering/selection made here will propagate to the next step

```
: ###############################################################
# display database with qgrid - this allows for interactive filtering
grid = qgrid.QGridWidget(df=df_criteria)
display(grid)
```

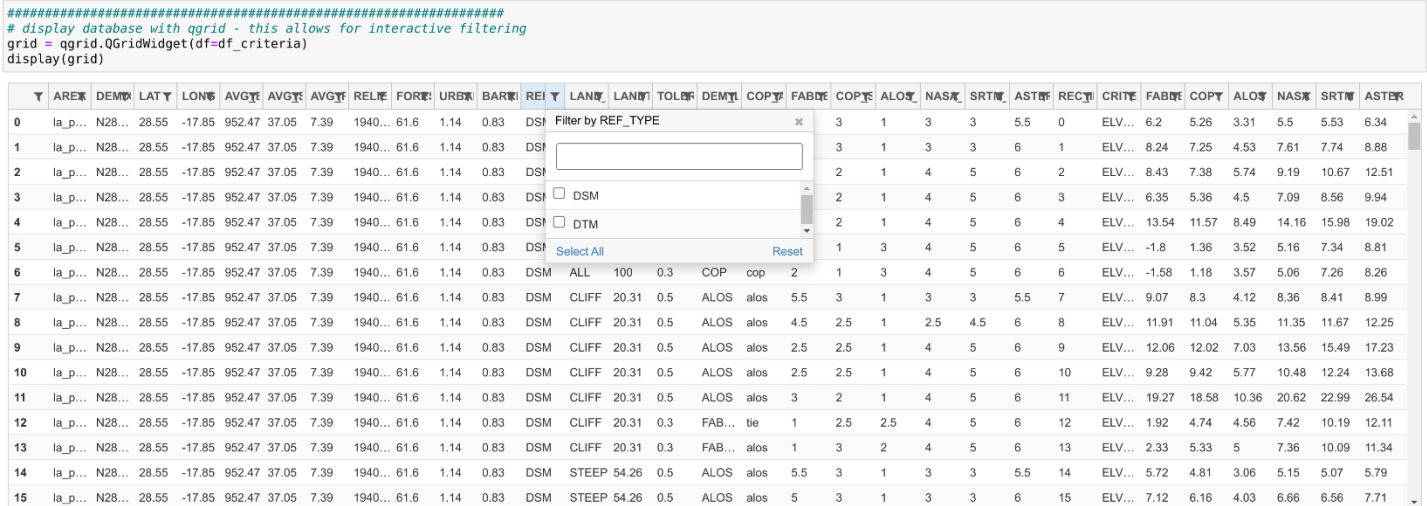| | AREA | DEM | LAT | LONG | AVG | AVG | AVG | RELI | FORE | URB | BARE | REF | LAND | LAND | TOLE | DEM | COP | FABD | COP | ALOS | NASA | SRTM | ASTE | RECT | CRITE | FABD | COP | ALOS | NASA | SRTM | ASTER |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | la_p... | N28... | 28.55 | -17.85 | 952.47 | 37.05 | 7.39 | 1940... | 61.6 | 1.14 | 0.83 | DSM | Filter by REF_TYPE | | | | | | 3 | 1 | 3 | 3 | 5.5 | 0 | ELV... | 6.2 | 5.26 | 3.31 | 5.5 | 5.53 | 6.34 |
| 1 | la_p... | N28... | 28.55 | -17.85 | 952.47 | 37.05 | 7.39 | 1940... | 61.6 | 1.14 | 0.83 | DSM | | | | ✕ | | | 3 | 1 | 3 | 3 | 6 | 1 | ELV... | 8.24 | 7.25 | 4.53 | 7.61 | 7.74 | 8.88 |
| 2 | la_p... | N28... | 28.55 | -17.85 | 952.47 | 37.05 | 7.39 | 1940... | 61.6 | 1.14 | 0.83 | DSM | | | | | | | 2 | 1 | 4 | 5 | 6 | 2 | ELV... | 8.43 | 7.38 | 5.74 | 9.19 | 10.67 | 12.51 |
| 3 | la_p... | N28... | 28.55 | -17.85 | 952.47 | 37.05 | 7.39 | 1940... | 61.6 | 1.14 | 0.83 | DSM | ☐ DSM | | | | | | 2 | 1 | 4 | 5 | 6 | 3 | ELV... | 6.35 | 5.36 | 4.5 | 7.09 | 8.56 | 9.94 |
| 4 | la_p... | N28... | 28.55 | -17.85 | 952.47 | 37.05 | 7.39 | 1940... | 61.6 | 1.14 | 0.83 | DSM | ☐ DTM | | | | | | 2 | 1 | 4 | 5 | 6 | 4 | ELV... | 13.54 | 11.57 | 8.49 | 14.16 | 15.98 | 19.02 |
| 5 | la_p... | N28... | 28.55 | -17.85 | 952.47 | 37.05 | 7.39 | 1940... | 61.6 | 1.14 | 0.83 | DSM | Select All | | Reset | | | | 1 | 3 | 4 | 5 | 6 | 5 | ELV... | -1.8 | 1.36 | 3.52 | 5.16 | 7.34 | 8.81 |
| 6 | la_p... | N28... | 28.55 | -17.85 | 952.47 | 37.05 | 7.39 | 1940... | 61.6 | 1.14 | 0.83 | DSM | ALL | 100 | 0.3 | COP | cop | 2 | 1 | 3 | 4 | 5 | 6 | 6 | ELV... | -1.58 | 1.18 | 3.57 | 5.06 | 7.26 | 8.26 |
| 7 | la_p... | N28... | 28.55 | -17.85 | 952.47 | 37.05 | 7.39 | 1940... | 61.6 | 1.14 | 0.83 | DSM | CLIFF | 20.31 | 0.5 | ALOS | alos | 5.5 | 3 | 1 | 3 | 3 | 5.5 | 7 | ELV... | 9.07 | 8.3 | 4.12 | 8.36 | 8.41 | 8.99 |
| 8 | la_p... | N28... | 28.55 | -17.85 | 952.47 | 37.05 | 7.39 | 1940... | 61.6 | 1.14 | 0.83 | DSM | CLIFF | 20.31 | 0.5 | ALOS | alos | 4.5 | 2.5 | 1 | 2.5 | 4.5 | 6 | 8 | ELV... | 11.91 | 11.04 | 5.35 | 11.35 | 11.67 | 12.25 |
| 9 | la_p... | N28... | 28.55 | -17.85 | 952.47 | 37.05 | 7.39 | 1940... | 61.6 | 1.14 | 0.83 | DSM | CLIFF | 20.31 | 0.5 | ALOS | alos | 2.5 | 2.5 | 1 | 4 | 5 | 6 | 9 | ELV... | 12.06 | 12.02 | 7.03 | 13.56 | 15.49 | 17.23 |
| 10 | la_p... | N28... | 28.55 | -17.85 | 952.47 | 37.05 | 7.39 | 1940... | 61.6 | 1.14 | 0.83 | DSM | CLIFF | 20.31 | 0.5 | ALOS | alos | 2.5 | 2.5 | 1 | 4 | 5 | 6 | 10 | ELV... | 9.28 | 9.42 | 5.77 | 10.48 | 12.24 | 13.68 |
| 11 | la_p... | N28... | 28.55 | -17.85 | 952.47 | 37.05 | 7.39 | 1940... | 61.6 | 1.14 | 0.83 | DSM | CLIFF | 20.31 | 0.5 | ALOS | alos | 3 | 2 | 1 | 4 | 5 | 6 | 11 | ELV... | 19.27 | 18.58 | 10.36 | 20.62 | 22.99 | 26.54 |
| 12 | la_p... | N28... | 28.55 | -17.85 | 952.47 | 37.05 | 7.39 | 1940... | 61.6 | 1.14 | 0.83 | DSM | CLIFF | 20.31 | 0.3 | FAB... | tie | 1 | 2.5 | 2.5 | 4 | 5 | 6 | 12 | ELV... | 1.92 | 4.74 | 4.56 | 7.42 | 10.19 | 12.11 |
| 13 | la_p... | N28... | 28.55 | -17.85 | 952.47 | 37.05 | 7.39 | 1940... | 61.6 | 1.14 | 0.83 | DSM | CLIFF | 20.31 | 0.3 | FAB... | alos | 1 | 3 | 2 | 4 | 5 | 6 | 13 | ELV... | 2.33 | 5.33 | 5 | 7.36 | 10.09 | 11.34 |
| 14 | la_p... | N28... | 28.55 | -17.85 | 952.47 | 37.05 | 7.39 | 1940... | 61.6 | 1.14 | 0.83 | DSM | STEEP | 54.26 | 0.5 | ALOS | alos | 5.5 | 3 | 1 | 3 | 3 | 5.5 | 14 | ELV... | 5.72 | 4.81 | 3.06 | 5.15 | 5.07 | 5.79 |
| 15 | la_p... | N28... | 28.55 | -17.85 | 952.47 | 37.05 | 7.39 | 1940... | 61.6 | 1.14 | 0.83 | DSM | STEEP | 54.26 | 0.5 | ALOS | alos | 5 | 3 | 1 | 3 | 3 | 6 | 15 | ELV... | 7.12 | 6.16 | 4.03 | 6.66 | 6.56 | 7.71 |

Figure 1.   GIS database rendered as a spreadsheet by qgrid. The selection between DTM and DSM is shown.

A pandas dataframe will be created from the GIS database, and the qgrid library will provide an interactive spreadsheet widget (Fig.1). This widget allows the user to query the database by various properties, such as selecting between DTM and DSM, filtering by DEMIX tile, area or difference metric.

It is imperative to select at least between DTM and DSM, as the analysis will not provide meaningful results if these two surface references are mixed.

A new dataframe is created based on the selection in the spreadsheet widget. This new dataframe will then be passed to the function responsible for the Friedman Test. This function will return information about the tolerances, which filters were applied to the database and the initial result, stating if the DEMs have statistically significant differences among them, and if the user can proceed to the post-hoc analysis (Fig.2).

```
Ranking with user-defined tolerances (might take a while...)

Filter settings for column REF_TYPE:['DSM']
Filter settings for column LAND_TYPE:['CLIFF', 'STEEP']

Results of the DEMIX Wine Contest

For k=6, CL=0.05, and N=1770, the critical value to compare is chi_crit=11.038
And since chi_r (5965.160) is greater than chi_crit (11.038)...
Yay!! We can reject the null hipothesis and go to the Post-Hoc analysis!!
```

Figure 2.   Example of results of the Friedman Test. Filters include a selection of DTM as surface type, CLIFF and STEEP as land types.

The next step is the post-hoc analysis, and the function will return a table (Fig.3), where each row corresponds to one DEM, and the columns are:

- rank – the final ranking, where lower is better;
- sum of ranks – numerical value of the sum of all ranks for the DEM;
- sum of ranks divided by number of opinions – this value might be used for comparing ranks made with different sets of criteria;

| | Rank | Sum of ranks | Sum of ranks divided by number of opinions | Ties with |
|---|---|---|---|---|
| **FABDEM** | 3.0 | 5958.0 | 3.366 | |
| **COP** | 2.0 | 4471.0 | 2.526 | |
| **ALOS** | 1.0 | 2372.5 | 1.340 | |
| **NASA** | 5.0 | 7194.0 | 4.064 | 5.0 |
| **SRTM** | 4.0 | 7182.0 | 4.058 | 4.0 |
| **ASTER** | 6.0 | 9992.5 | 5.645 | |

- ties with – with which DEM there is a 'tie', that is, these DEMs don't have a statistically significant difference

Figure 3.   Example of the post-hoc analysis results, using the same filters as in Fig.3. Here SRTM and ASTER GDEM are tied in 4[th] place.
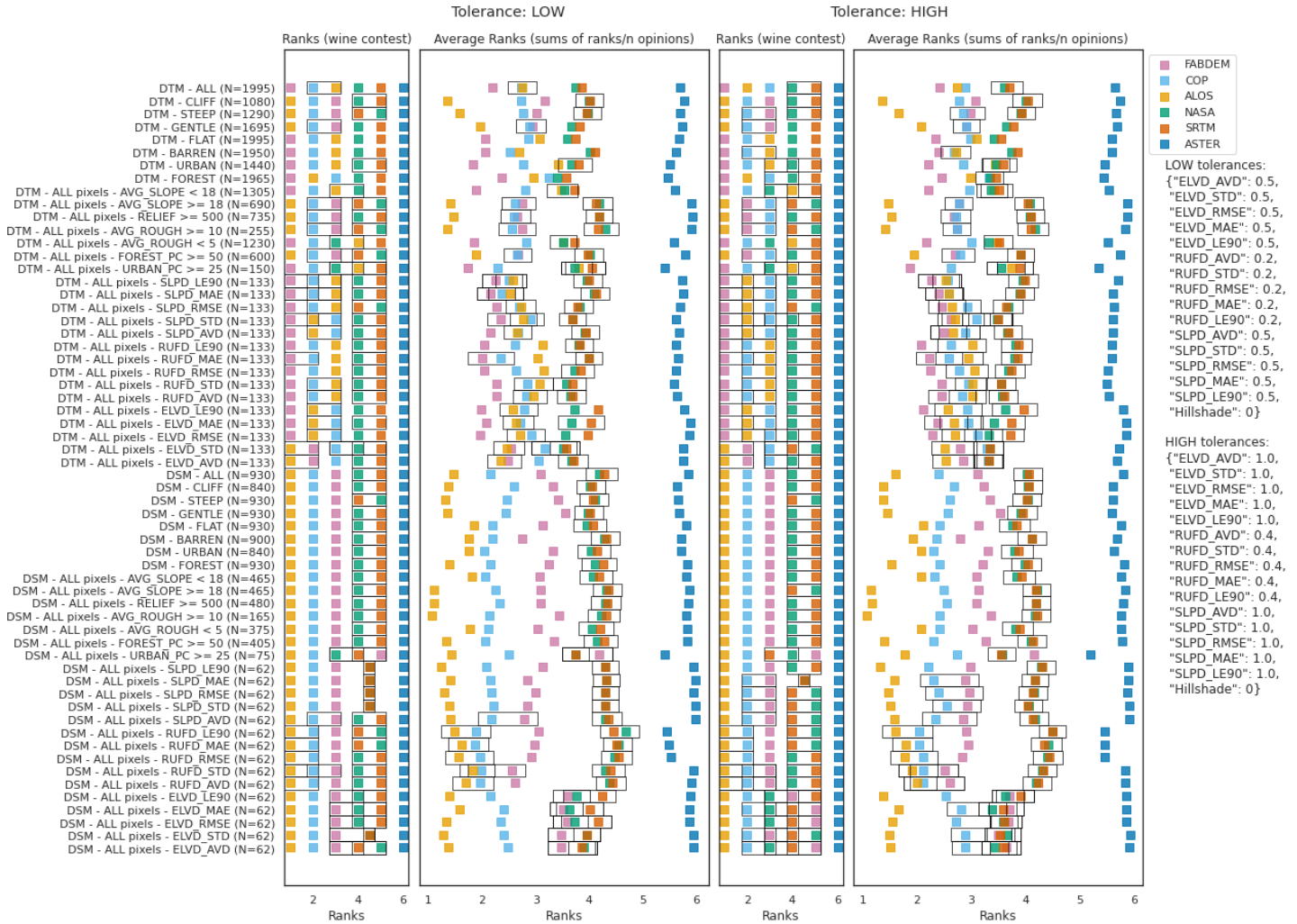
Figure 4. Scatterplot of 60 Wine Contests ran with different sets of criteria and tolerances. Rectangles over DEMs' symbols indicate a tie between them.

The notebook also provides tools to analyze outputs by creating graphics and figures. The plots included in the notebook were intended to illustrate the publication [7], and the choice of predefined parameters of colors, symbols, and text annotations reflect this.

Included in the notebook are a plot of correlation matrices, a customized scatterplot (Fig.4) of 60 wine contest ranks (including database filters on surface type, land type, land cover, geomorphometric indices, and two sets of tolerances), scatterplots of signed metrics (means, medians), unsigned metrics (RMSE, standard deviation), and of selected criteria per selected tiles (chosen as representative of the metrics' behavior). All figures can be saved in common formats (e.g., .png, .svg). Note that the plots are not produced in 'publication-ready' formatting, as the author prefers to finalize the figures' in an illustration software.

## IV. CONCLUSION

This paper presented the DEMIX Wine Contest Jupyter Notebook, an open-source tool developed to provide an interface to explore the DEMIX GIS database and to generate Wine Contest results from several sets of criteria and tolerances. The notebook is available in GitHub. Indications of errors, bugs or suggestions to improve the code are welcome.

## V. ACKNOWLEDGEMENTS

## REFERENCES

[1]  Shen, H. Interactive notebooks: Sharing the code. 2014. *Nature* **515**, 151–152. https://doi.org/10.1038/515151a

[2]  Rule A, Birmingham A, Zuniga C, Altintas I, Huang S-C, Knight R, et al. (2019) Ten simple rules for writing and sharing computational analyses in Jupyter Notebooks. PLoS Comput Biol 15(7): e1007007. https://doi.org/10.1371/journal.pcbi.1007007

[3]  Kluyver T., Ragan-Kelley B., Pérez F., Granger B.E., Bussonnier M., Frederic J., Kelley K., Hamrick J.B., Grout J., Corlay S. et al. 2016. Jupyter notebooks - a publishing format for reproducible computational workflows. In: Loizides F., Scmidt B. (eds) Positioning and Power in Academic Publishing: Players, Agents and Agenda. IOS Press, pp 87–90. http://dx.doi.org/10.3233/978-1-61499-649-1-87

[4]  Knuth, D.E. 1984. Literate Programming, *The Computer Journal*, 27:2, 97–111, https://doi.org/10.1093/comjnl/27.2.97

[5]  Pimentel, J.F., Murta, L., Braganholo, V. *et al.* 2021. Understanding and improving the quality and reproducibility of Jupyter notebooks. Empirical Software Engineering  26:65. https://doi.org/10.1007/s10664-021-09961-9

[6]  Strobl, P.A.; Bielski, C.; Guth, P.L.; Grohmann, C.H.; Muller, J.P.; López-Vázquez, C.; Gesch, D.B.; Amatulli, G.; Riazanoff, S.; Carabajal, C. 2021.The Digital Elevation Model Intercomparison eXperiment DEMIX, a community based approach at global DEM benchmarking. Int. Arch. Photogramm. Remote Sens. Spat. Inf. Sci. 2021, XLIII-B4-2021, 395–400. https://doi.org/10.5194/isprs-archives-XLIII-B4-2021-395-2021

[7]  Bielski, C.; López-Vázquez, C.; Grohmann, C.H.; Guth, P.L.; The TMSG DEMIX Working Group. 2023 (preprint). DEMIX Wine Contest Method Ranks ALOS AW3D30, COPDEM, and FABDEM as Top 1" Global DEMs. ArXiV. https://arxiv.org/abs/2302.08425

[8]  Grohmann, C.H. 2023. DEMIX Wine Contest Jupyter Notebook. URL: https://github.com/CarlosGrohmann/DEMIX_wine_contest

[9]  Python Software Foundation, 2021. Python Programming Language, version 3.10.x. http://www.python.org/

[10] Friedman, M. 1937. The Use of Ranks to Avoid the Assumption of Normality Implicit in the Analysis of Variance. Journal of the American Statistical Association 32, 675-701. https://doi.org/10.2307/2279372

[11] Dunn, O.J. 1961. Multiple Comparisons among Means. Jour. of the Am. Statistical Association 56,52-64. https://doi.org/10.2307/2282330

[12] qgrid python library. 2020. https://github.com/lukewys/qgrid

[13] Harris, C. R., Millman, K. J., van der Walt, S. J., Gommers, R., Virtanen, P., Cournapeau, D., Wieser, E., Taylor, J., Berg, S., Smith, N. J., Kern, R., Picus, M., Hoyer, S., van Kerkwijk, M. H., Brett, M., Haldane, A., del Río, J. F., Wiebe, M., Peterson, P., Gérard-Marchant, P., Sheppard, K., Reddy, T., Weckesser, W., Abbasi, H., Gohlke, C., Oliphant, T. E., 2020. Array programming with numpy. Nature 585 (7825), 357–362. https://doi.org/10.1038/s41586-020-2649-2

[14] McKinney, W., 2011. pandas: a Foundational Python Library for Data Analysis and Statistics. In: Python for High Performance and Scientific Computing. https://www.dlr.de/sc/portaldata/15/resources/dokumente/pyhpc2011/submissions/pyhpc2011_submission_9.pdf

[15] Hunter, J.D., 2007. Matplotlib: A 2D graphics environment. Computing. Science & Engineering 9(3),90-95. https://doi.org/10.1109/MCSE.2007.55

[16] Waskom, M.L., 2021. seaborn: statistical data visualization. Journal of Open Source Software, 6(60), 3021, https://doi.org/10.21105/joss.03021

[17] Guth, P.L. 2022. DEMIX GIS Database (1.0) [Data set]. Zenodo. https://doi.org/10.5281/zenodo.7402618

[18] Guth, P.L., 2023. MICRODEM. https://github.com/prof-pguth/git_microdem