

Welcome to 0123-import-scripts!

This is a template repository that can be used for the archiving process of a big dataset at the end of a research project's lifetime. Download and unpack this repository from the [excel2xml documentation page](#).

In this README, you will learn how to write a Python script for preparing data for an import into DSP.

Featuring:

- first steps with Visual Studio Code
- the module `excel2xml` of dsp-tools
- the benefits of Version Control with Git
- the benefits of the Debugging Mode
- extras: OpenRefine, Git GUIs, regexpr

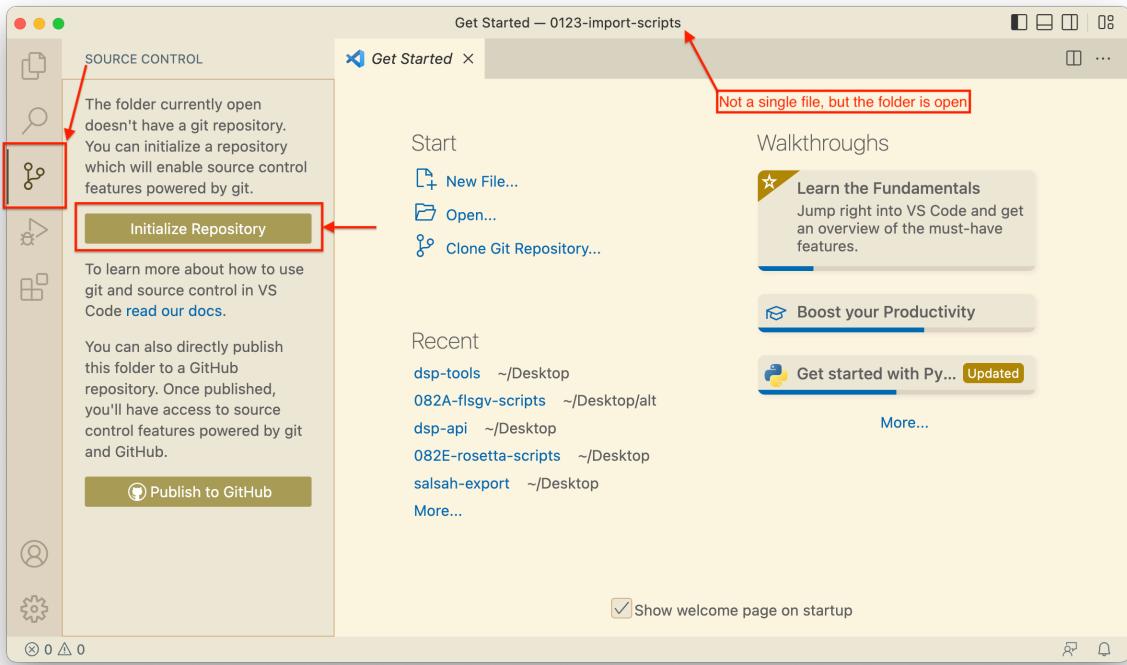
First steps with Visual Studio Code (VSC)

[Visual Studio Code](#) is the industry standard IDE, free to use, and recommended for your daily work at DaSCH. Recommended extensions to install from within VSC:

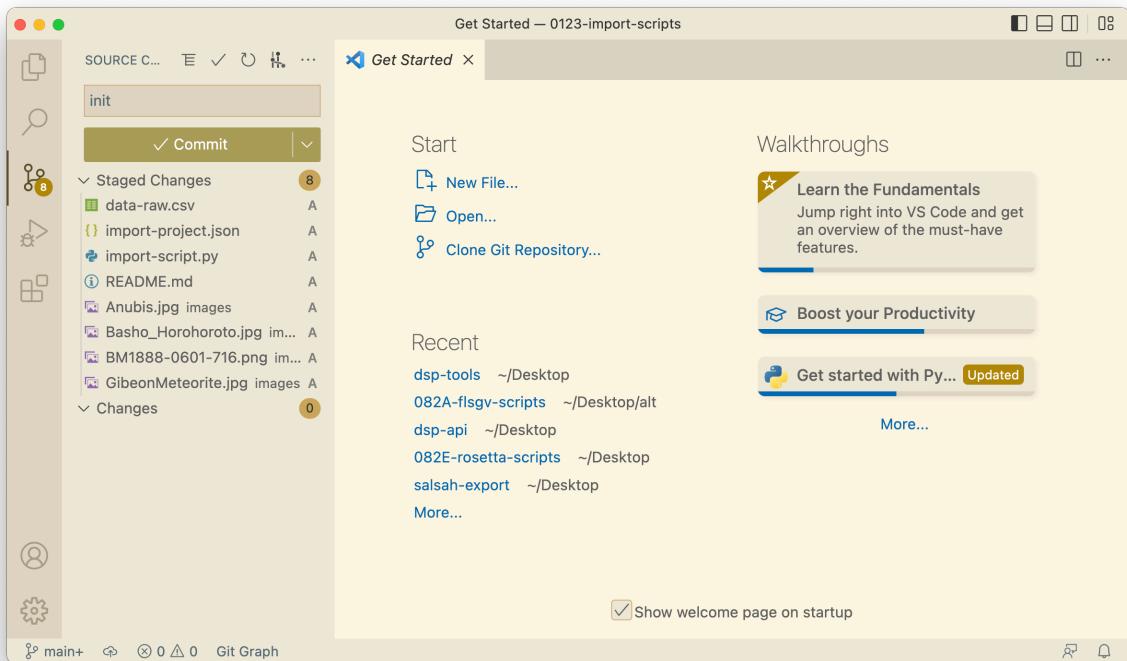
- redhat.vscode-xml
- ms-python.python
- ms-python.vscode-pylance
- zainchen.json
- nickdemayo.vscode-json-editor
- visualstudioexptteam.vscodeintellicode
- redhat.vscode-xml
- dotjoshjohnson.xml

Initialize Git

Open this repository in Visual Studio Code, change to the "Source Control" tab, and click on "Initialize Repository":



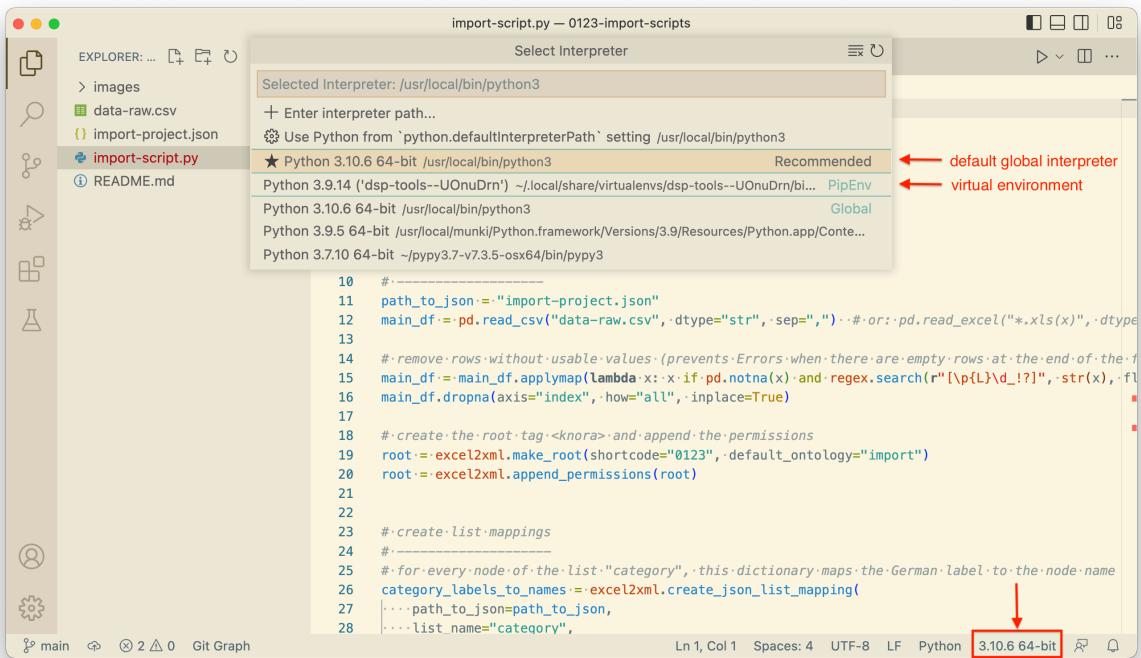
Stage all changes, write "init" as commit message, and commit all changes:



You now have the option "Publish Branch". This is to synchronize your local work with a GitHub repository on <https://github.com/dasch-swiss/>. For this purpose, replace `0123` by your project's shortcode, and `import` by your project's shortname. This is especially recommended for big projects where you spend weeks/months on, when you might want to have a backup, or when you want to invite colleagues for collaboration or a code review.

Choose a Python interpreter

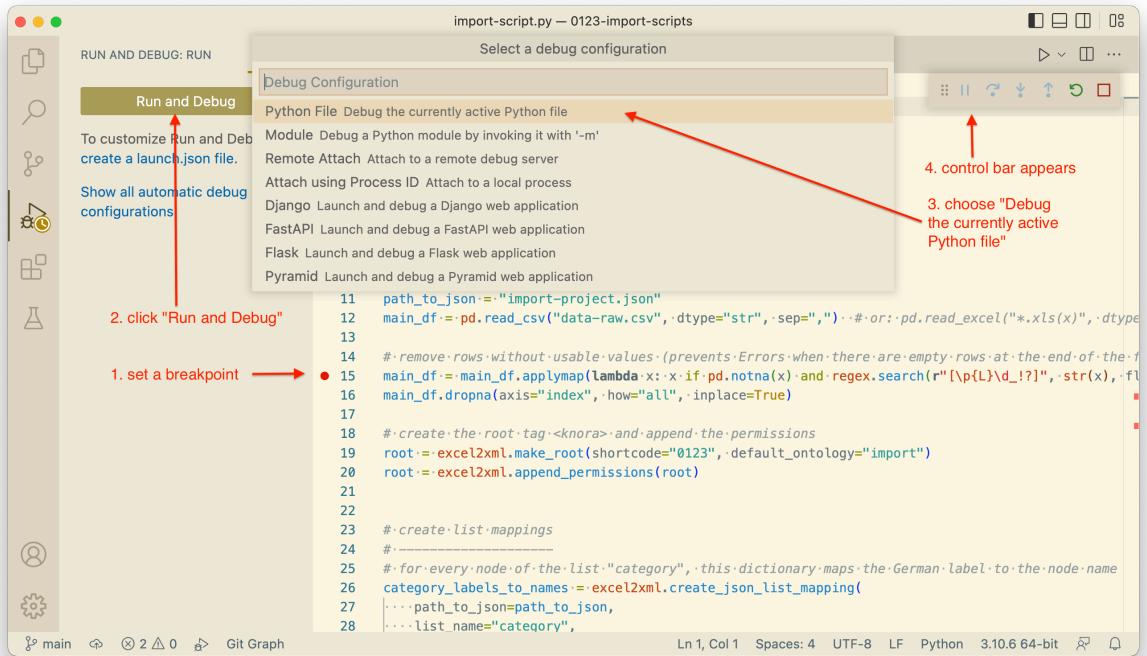
Open `import-script.py`. You can now choose a Python interpreter by clicking on the Version number on the bottom right. You can either work with the global (system-wide) Python, or you can create a [virtual environment](#) for your project. If you don't know which one to choose, take the one installed via Homebrew, which is located in `/usr/local/Cellar`. Probably you already have a symlink (`/usr/local/bin/python3` or `/usr/bin/python3`) that redirects to `/usr/local/Cellar`. The only thing that you shouldn't do is selecting a virtual environment of another project.



The benefits of the debugging mode

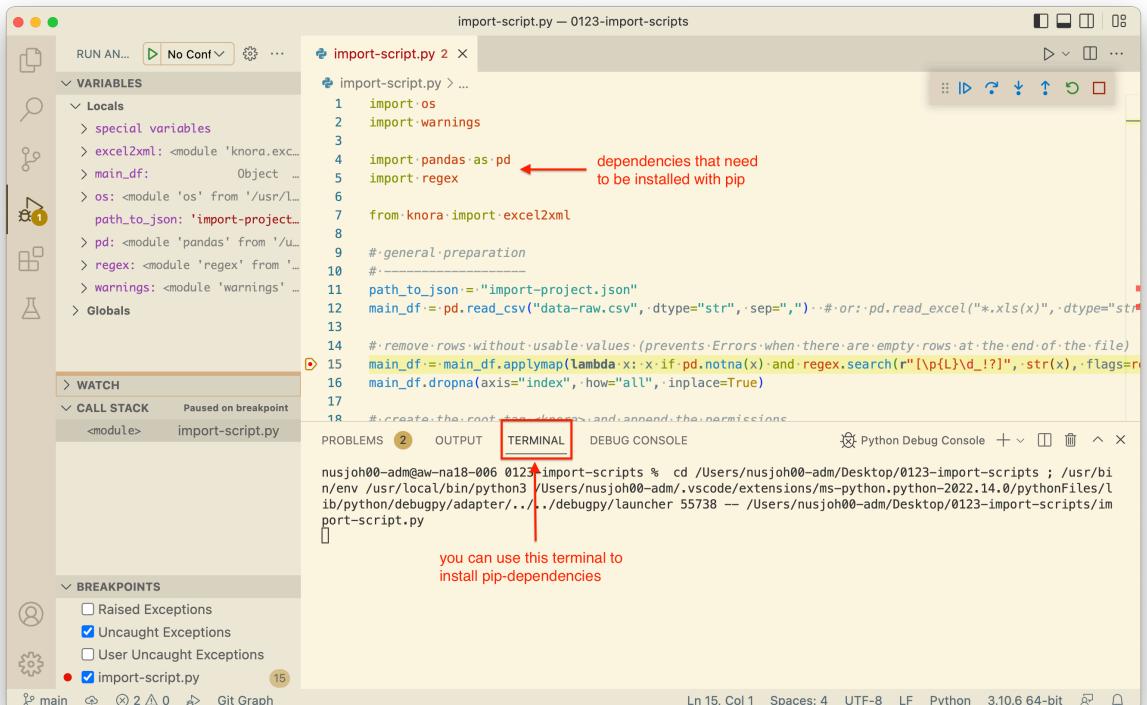
To start the debugging process, switch to the "Run and Debug" tab.

1. set a break point
2. click "Run and Debug"
3. choose "Debug the currently active Python file"
4. The control bar appears, and debugging starts.



Code execution will interrupt at your break point, that means, before the line of the break point is executed. Use this opportunity to inspect what has been done until now in the "Variables" area on the left, where the current state of the program is shown.

If one of the dependencies is not installed, install it with `pip install package` in the Terminal of Visual Studio Code.



If you want to experiment with different scenarios how to proceed, go to the "Debug Console" where you can execute code. For example, let's inspect the Pandas Dataframe by typing `main_df.info`.

You see that there are some empty rows at the end which don't contain useful data. The next two lines of code will eliminate them. Click on "Step Over" two times, or set a new break point two lines further down and click on "Continue". Now, type again `main_df.info` in the Debug Console. You will see that the empty rows are gone.

```
import os
import warnings
import pandas as pd
import regex
from knora import excel2xml
# general preparation
path_to_json = "import-project.json"
main_df = pd.read_csv("data-raw.csv", dtype="str", sep=",") # or: pd.read_excel("*.xls(x)", dtype="str")
# remove rows without usable values (prevents Errors when there are empty rows at the end of the file)
main_df = main_df.applymap(lambda x: x if pd.notna(x) and regex.search(r"\[p(L)\d_!\]", str(x), flags=regex.U) else pd.NA)
main_df.dropna(axis="index", how="all", inplace=True)
# create the root tag <knora> and append the permissions
root = excel2xml.make_root(shortcode="0123", default_ontology="import")
```

	Object	...	NaN	NaN	
0	Anubis	Bengal cat	...	NaN	NaN
1	Meteorite	Gibeon Meteorite	...	NaN	NaN
2	BM1888-0601-716	Lekythos	...	NaN	NaN
> 3	Horohoro	Picture and Poem by Matsuo Bashō	...	NaN	NaN
4			NaN	...	NaN
5			NaN	...	NaN

You see that the debugging mode is a useful tool to understand code and to inspect it for correctness.

Tip

Make regular use of the debugging mode to check if your code really does what you think it should do!

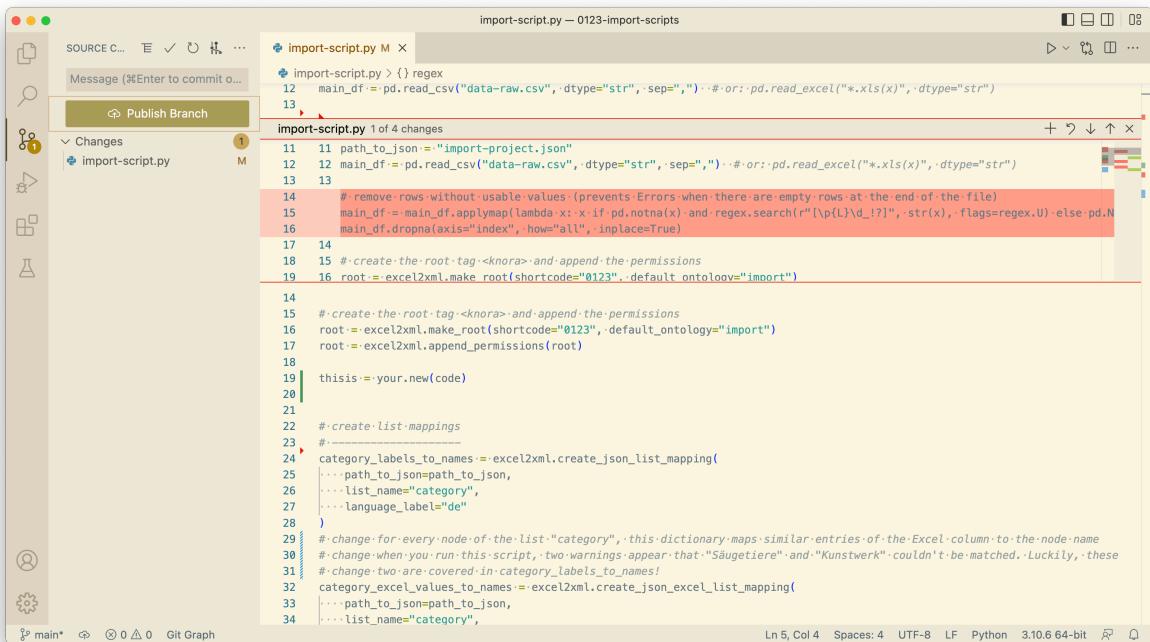
The benefits of version control

One of the big benefits of version control is the diff viewer. Visual Studio highlights the changes you have introduced since your last commit.

- Deletions are shown as red triangle.
- Additions are shown as green bars.
- Changed lines are shown as striped bars.

Click on these visual elements to see a small popup that shows you the difference. In the

popup, you can stage the change, revert it, or jump to the next/previous change.



The screenshot shows a GitHub commit interface. On the left, there's a sidebar with icons for file operations like copy, move, delete, and merge. The main area shows a file named 'import-script.py' with 14 changes staged. The code is as follows:

```
import-project.json"
11 main_df = pd.read_csv("data-raw.csv", dtype="str", sep=",") # or: pd.read_excel("*.xls(x)", dtype="str")
12 main_df = pd.read_csv("data-raw.csv", dtype="str", sep=",") # or: pd.read_excel("*.xls(x)", dtype="str")
13
14 # remove rows without usable values (prevents Errors when there are empty rows at the end of the file)
15 main_df = main_df.applymap(lambda x: x if pd.notna(x) and not regex.search(r"\p{L}\d_?!", str(x), flags=regexec.U) else pd.N
16 main_df.dropna(axis="index", how="all", inplace=True)
17
18 # create the root tag <knora> and append the permissions
19 root = excel2xml.make_root(shortcode="0123", default_ontology="import")
20
21
22 # create list mappings
23
24 category_labels_to_names = excel2xml.create_json_list_mapping(
25     ...path_to_json=...path_to_json,
26     ...list_name="category",
27     ...language_label="de"
28 )
29 # change for every node of the list "category", this dictionary maps similar entries of the Excel column to the node name
30 # change when you run this script, two warnings appear that "Säugetiere" and "Kunstwerk" couldn't be matched. Luckily, these
31 # change two are covered in category_labels_to_names!
32 category_excel_values_to_names = excel2xml.create_json_excel_list_mapping(
33     ...path_to_json=...path_to_json,
34     ...list_name="category",
```

At the bottom, status indicators show 5 lines, 4 columns, spaces, and encoding (UTF-8). A note says 'Ln 5, Col 4'.

Once you have a bunch of code changes that can be meaningfully grouped together, you should make a commit (and perhaps push it to a GitHub repo).

Tips

Test your code (e.g. with the debugging mode) before committing it.

Make small commits that contain only one new feature.

Some extras

Data cleaning with OpenRefine

[OpenRefine](#) is a tool for working with messy data. Once downloaded and installed, it runs as a local server, accessed by your browser. So, all data remains on your own machine.

Installation is quick and painless: `brew install openrefine`

The potentials for the everyday work of the Client Services at DaSCH are twofold: 1. Data cleaning (recommended): For this purpose, you can think of OpenRefine as a much better version of Excel. You can perform operations which would be very tiresome in Excel. 2. Conversion to our dps-customised xml format for bulk upload (not recommended)

Read more in [this report](#).

Git GUIs

Git can be complicated, so you will appreciate to work with one of these GUIs:

- [SmartGit](#) has a free edition.
- [GitHub Desktop](#)
- [SourceTree](#)

Learn, build and test RegEx

<https://regexr.com>