# Data mining versus manual screening to select papers for inclusion in systematic reviews: a novel method to increase efficiency

Elena Ierardi, J Chris Eilbeck, Frederike van Wijck, Myzoon Ali, Fiona Coupar

# Instructions for Windows

## Use of the Command Line

You will need to know where the command line is and how to type commands here (Figure 1). Click on the Start button in Windows and type "**cmd"** in the search option at the bottom where it says "Type here to search" and press Enter.
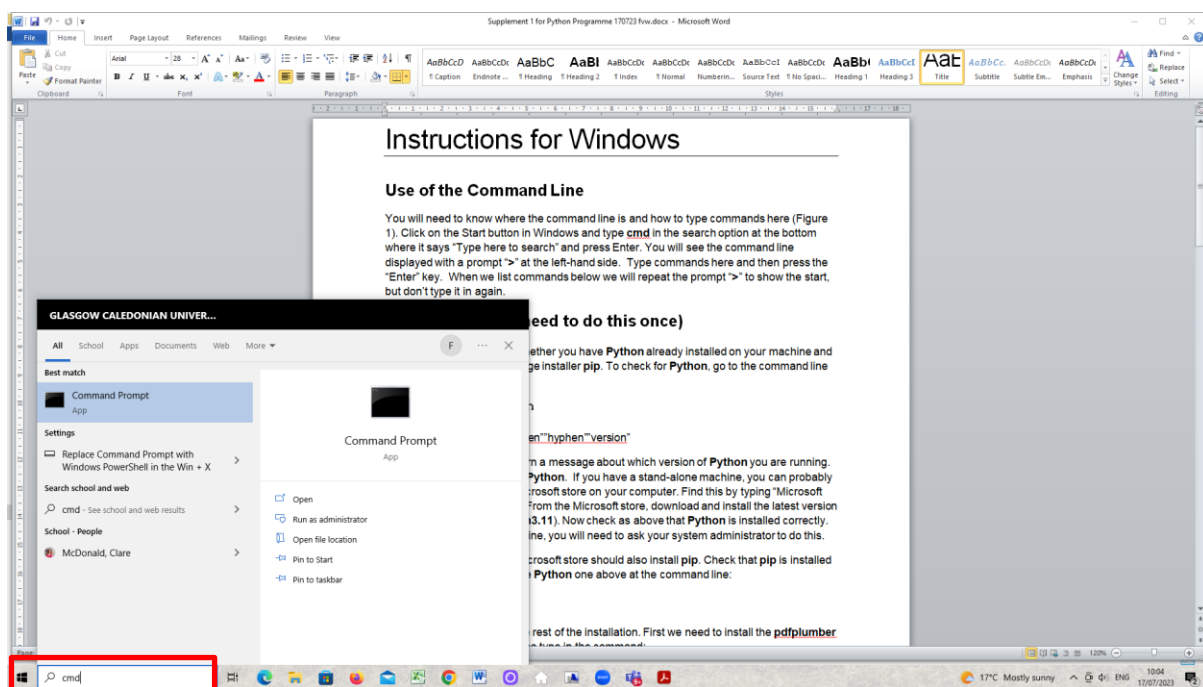


**Figure 1. Typing "cmd" into the 'Type here to search' window**

You will see the command line displayed with a prompt "**>**" at the left-hand side (Fig. 2). Type commands here and then press the "Enter" key.  When we list commands below we will repeat the prompt "**>**" to show the start, but don't type it in again.
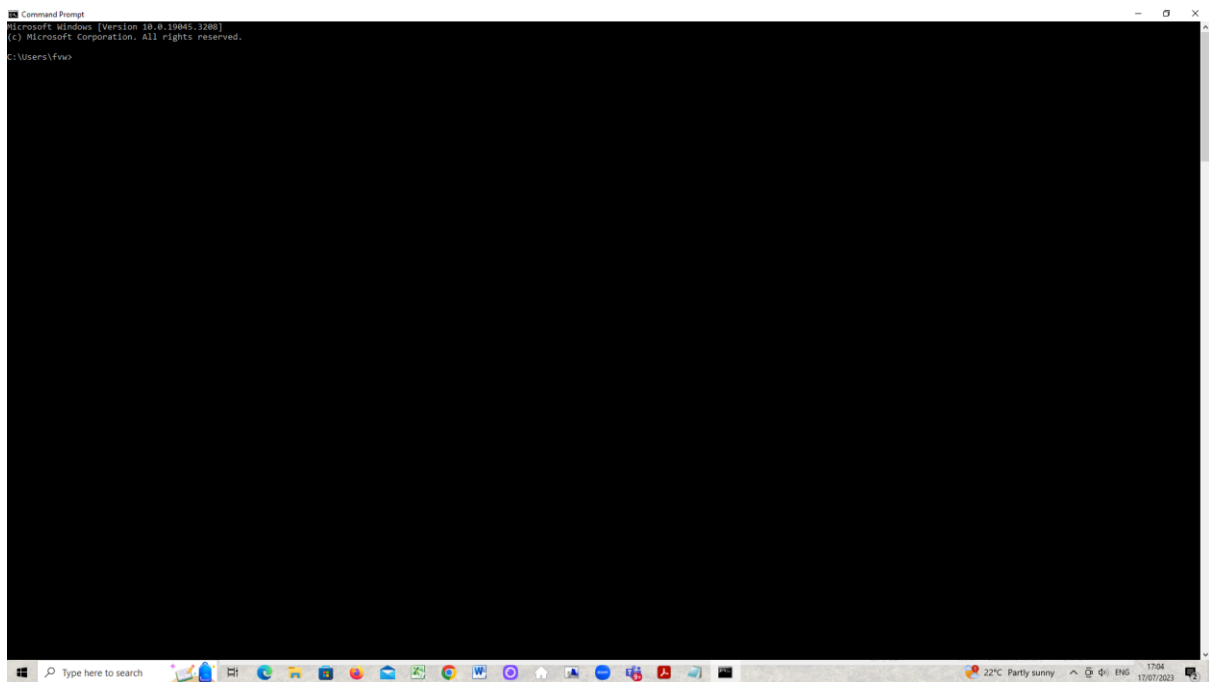


**Figure 2. Display of the command line in Command Prompt**

# **Preparation** (you only need to do this once)

The next steps depend on whether you have **Python** already installed on your machine and in addition the **Python** package installer **pip**.

## **Checking Python**

To check for **Python**, go to the command line and type:

```
Python --version
```

This is "Python""space""hyphen""hyphen""version"

If successful, this should return a message about which version of **Python** you are running (Fig. 3).
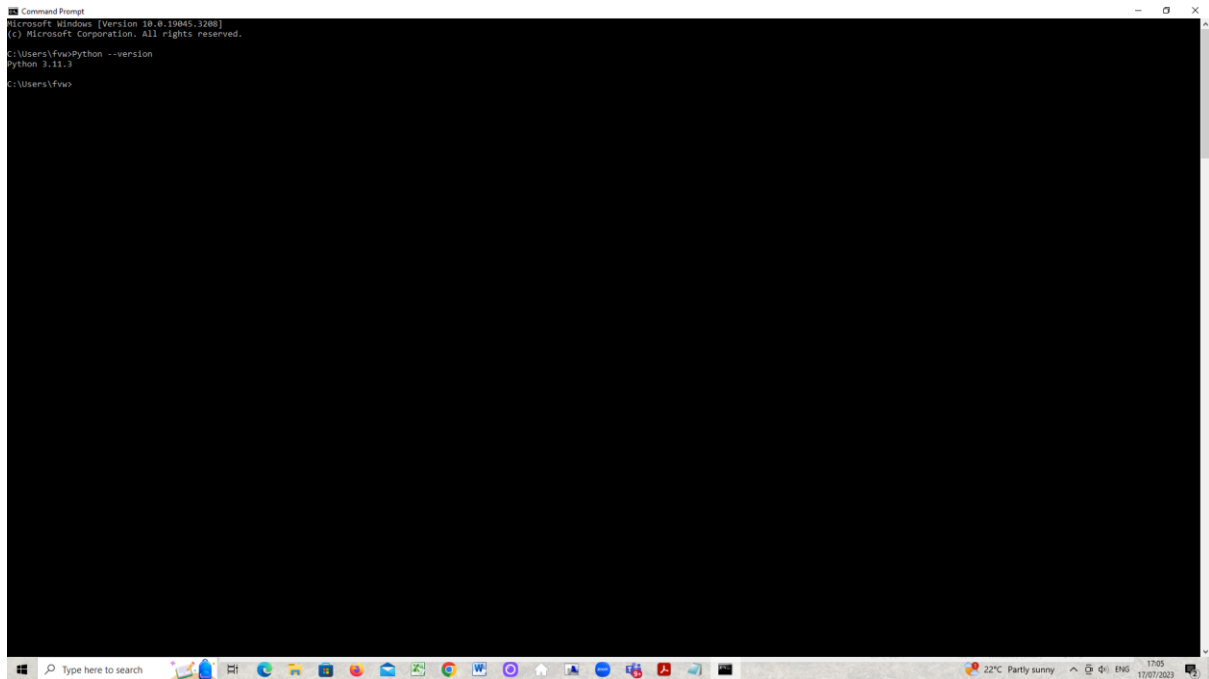
**Figure 3. Response indicating the available version of Python**

# Installing Python

If you don't have **Python** installed yet, you will need to do so.  If you have a stand-alone machine, you can probably do this yourself. Go to the Microsoft store on your computer. Find this by typing "MS store" in the search window (Fig. 4).
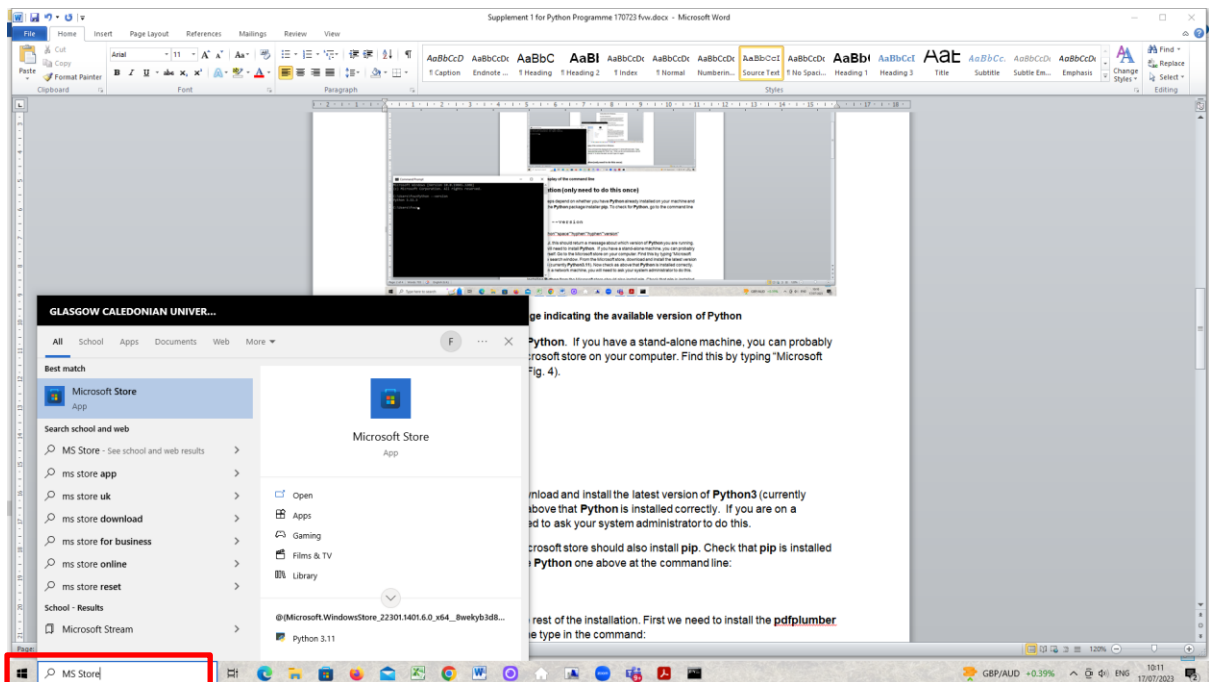


**Fig. 4. Typing "MS Store" into the search window**

From the Microsoft store, download and install the latest version of **Python3** (currently **Python3.11**). Now check as above that **Python** is installed correctly.  If you are on a network machine, you may need to ask your system administrator to do this.

Installing **Python** from the Microsoft store should also install **pip**. Check that **pip** is installed with a command similar to the **Python** one above at the command line:

```
pip --version
```

We can now proceed with the rest of the installation. First we need to install the **pdfplumber** package.  At the command line type in the command:
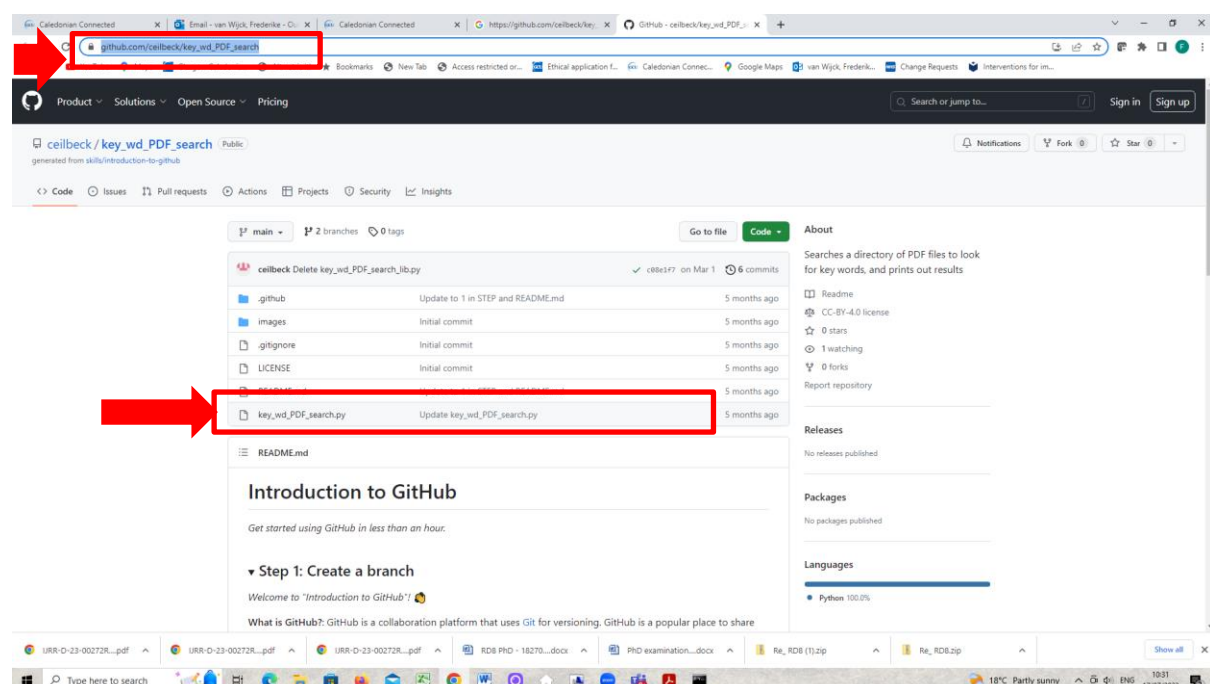
```
> pip install pdfplumber
```

It may take a little while, but you should get eventually a message saying "Successfully installed **pdfplumber**".

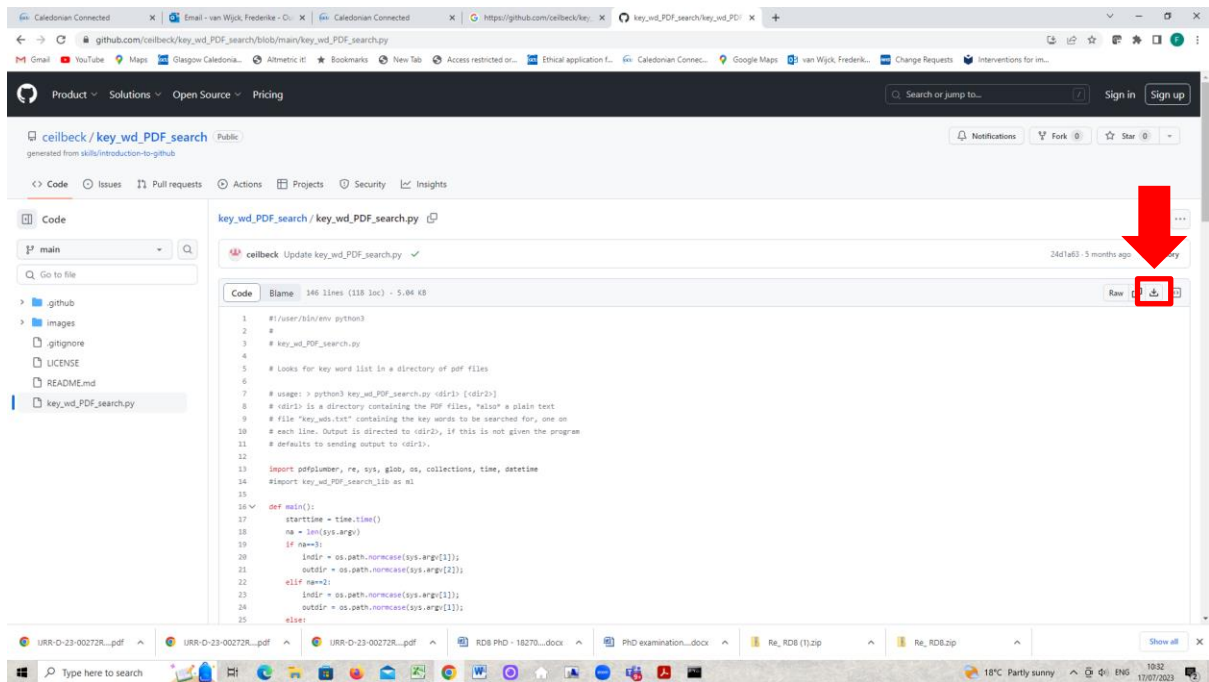## Downloading the Python programme to search for key words

Copy this address into your web browser (Fig. 5, top arrow):

https://github.com/ceilbeck/key_wd_PDF_search



**Fig. 5. Searching for the Python programme to search for key words on Github**

Then click on the hyperlink `key_wd_PDF_search.py`  (Fig. 5, bottom arrow). This takes you to the actual programme.

4

**Fig. 6. Finding the Python programme to search for key words**

Now click on the 'Download' icon, which downloads the programme (`Fig. 6`). This will need to be stored in a specific directory, as explained in the next step.

## Creating a directory for your work

Now you need to set up a suitable directory system for you work. Still on the command line, enter:

```
cd ..
```

That's **cd** "space" followed by two dots. This should take you to your main directory on your **C** drive, something like **c:\\user >**

At the command line, create a subfolder and name it "**programs**":

```
mkdir programs
```

move into this subfolder:

```
cd programs
```

and create another subfolder:

```
mkdir PDFfiles
```

## Moving files into their correct folders

It is crucial that the various files are in their designated folders, otherwise the program is unable to find the necessary information to do its work:

1. Copy the following files in the "**programs**" folder you created earlier:
   1. this Instructions file (from this article's Supplementary files) and
   2. the **key_wd_PDF_search.py** file (downloaded from Github, as explained above)
2. Copy the **key_wrd.txt** in the "**PDFfiles**" subfolder (This is the sub folder of the "**programs**" folder).
3. Change the keywords as per your requirements. To do that, open the **key_wrd.txt** from the "**PDFfiles**" subfolder, amend any of the keywords and save the file (make sure to retain the file name and its extension (i.e. ".txt"), or the computer program won't work as it does not recognise the file.
4. Put all of your PDF files (on which you want to do a keyword search) in "**PDFfiles**" subfolder

# Running the data mining programme

You are now ready to run the data mining programme:

1. Open the **Command Prompt** in Windows (see above if you are not sure how) and type:

   ```
   cd ..
   ```
   (that's **cd** "space" followed by two dots)

   into the Command Prompt and press the **Enter** key on your keyboard (repeat if necessary until it displays **C:\>** to show that you are in the main C: drive).
2. Switch to the folder "**programs**" by typing:

   ```
   cd programs
   ```
   (that's **cd** "space" followed by "programs").

   into the Command Prompt and press the **Enter** key on your keyboard.

**3.** Now type:

   ```
   Python3 key_wd_PDF_search.py PDFfiles
   ```

   into the Command Prompt and press the **Enter** key on your keyboard.
4. The Data Mining Tool will now run and the screen will be blank during this time. When it completes (Fig. 7), the generated output text files (one for each paper processed) will be ready to view in the subfolder "**PDFfiles**" (Fig. 8).
5. A summary of the resulting can also be found in "Notes.txt" in the same folder.

*NOTE: Please note that files which are password-protected or that are images (scanned) are not going to be processed and these output files will be empty. The program may skip out of a PDF being processed if it encounters a problem, so not all pages may be returned.*

**Fig. 7. Screenshot of the Python data mining programme having completed**



**Fig. 8. Screenshot of the Output files and Notes files added to the PDF files subfolder**

Figure 9 shows what the output from a single PDF may look like, and Fig. 10 shows the contents of the Notes files after three PDFs have been processed; this file records the key words, the number of times they have been identified in all PDF files analysed, and processing time.
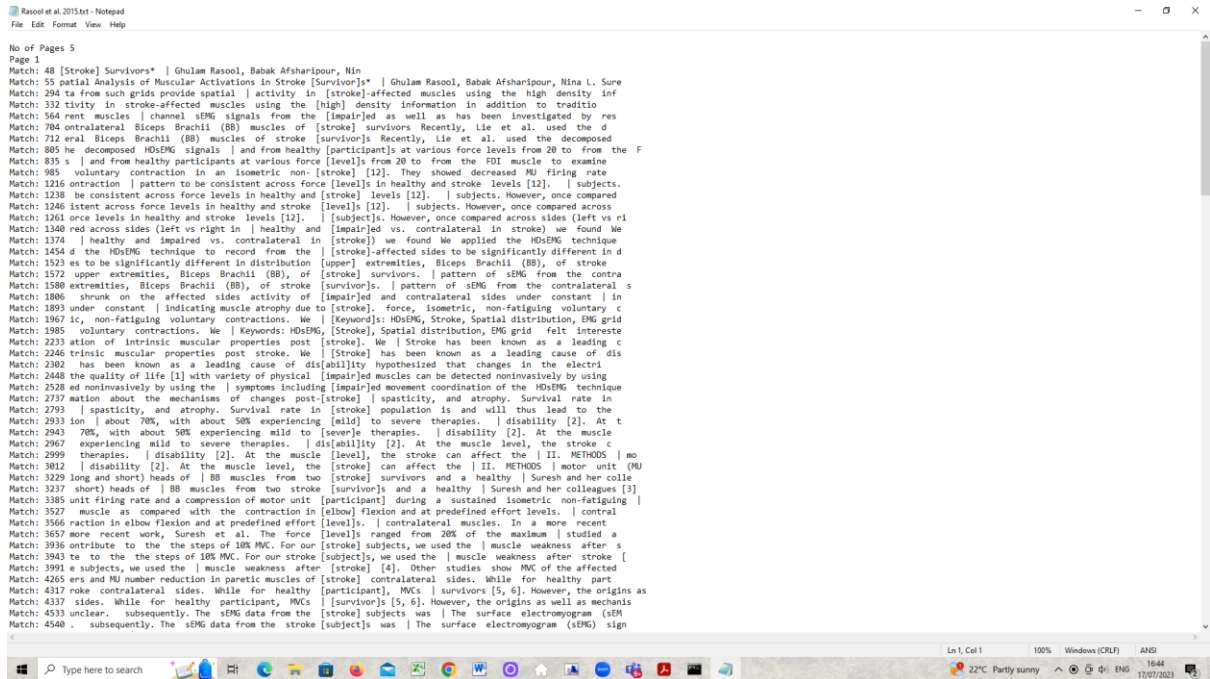
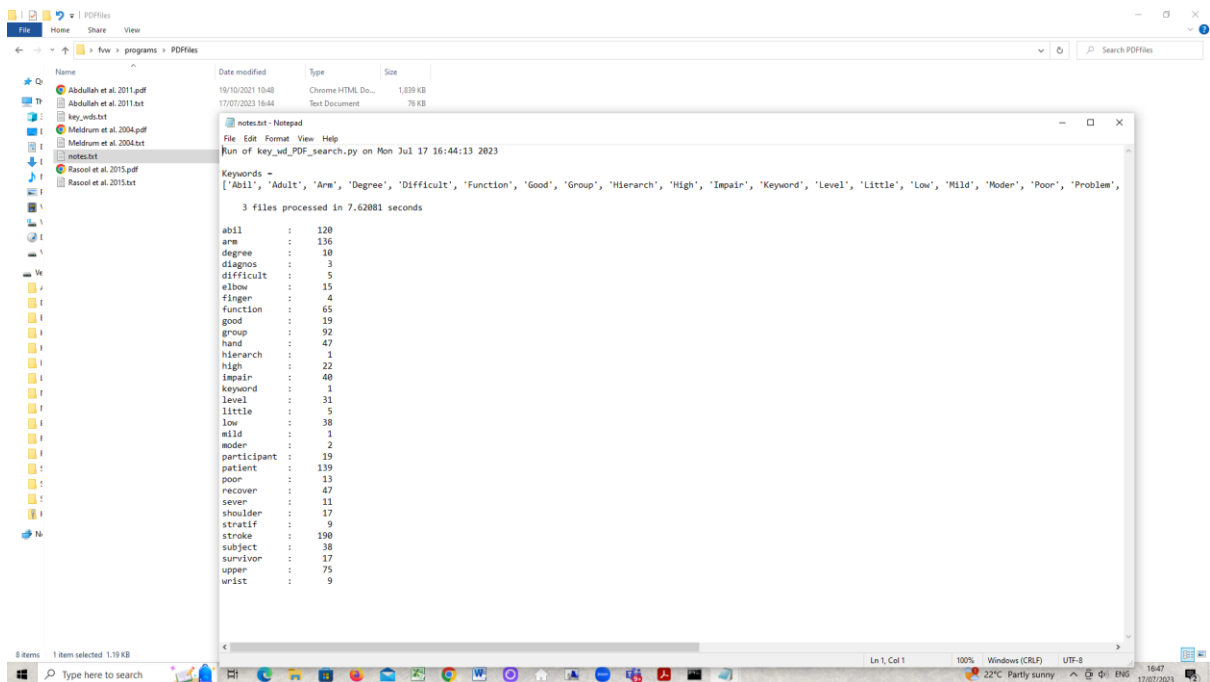**Fig. 9. Screenshot of an excerpt of one of the Output files**



**Fig. 10. Screenshot of the Notes file, based on a sample of three processed PDF papers.**