**A Note on the Files in the 01code Folder**

The folder contains the python code used to filter and analyse the raw database, as well as the published R Markdown document containing the code used to produce the plots.

**A Note on the Files in the 02Processed Folder**

This folder contains the initial dataset as produced by the python code, and the dataset in long format, which is the format used for Plot 1. These were uploaded separately to the R Markdown document to save time converting the data from the initial format to the long format. After some tests this seemed to save 30 seconds or so in loading the document.

**Converting the Caissabase file to utf-8**

The Caissabase dataset seems to require being converted from ISO8859-1 (AKA Latin1) encoding to utf-8 to work with the python code that filters and analyses it.

This can be done in a linux terminal with the following code:

1.
```
file -i file.name
```
[checks the encoding - returns "8859_1"]

2.
```
iconv -f 8859_1 -t UTF-8//TRANSLIT file.name -o file.name
```
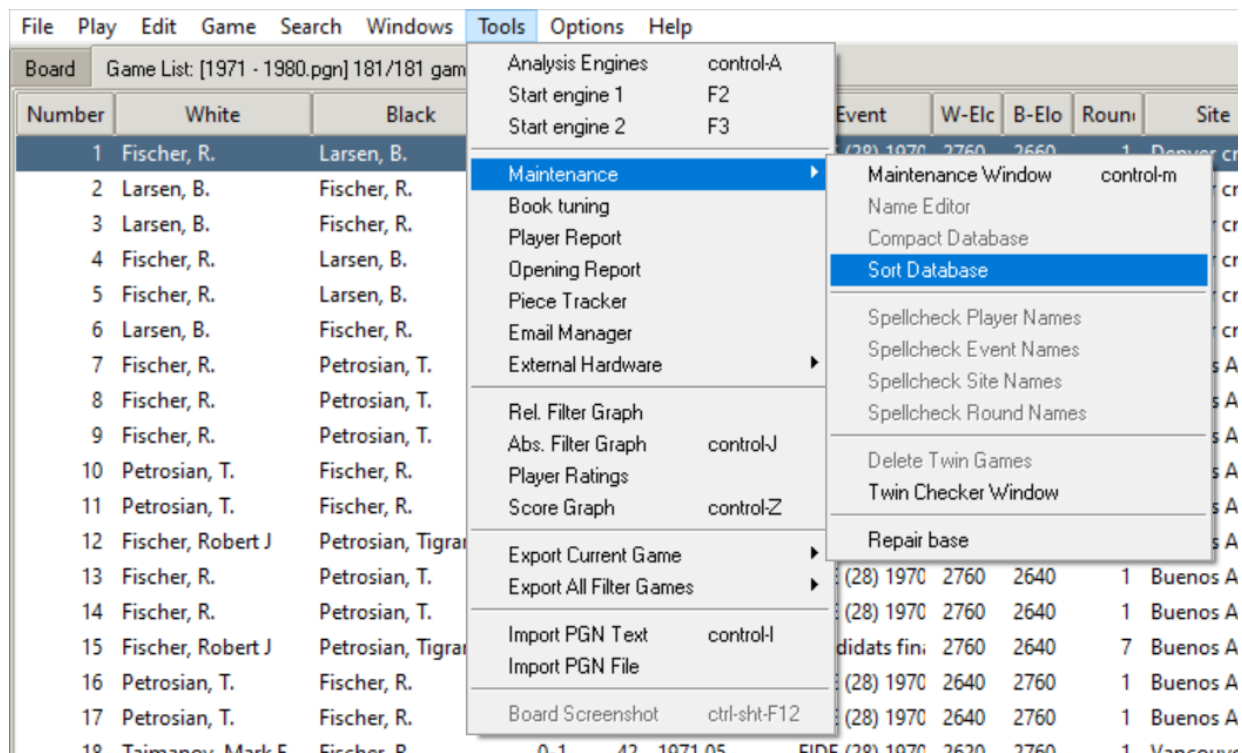
This can also be done inside python.

_____

**Alternative Sorting Methods**

Another method to sort the full Caissabase data, which depending on the purpose may be faster than using the included python code, is to download the Scid, a programme for interfacing with pgn files: http://scidvspc.sourceforge.net/

Upload the .pgn file, and go to Game -> List All Games.

The games can then be filtered according to a range of information using the 'Sort Database' window. Found in Tools -> Maintenance -> Sort Database

Scid allows you to remove multiple games at once, resort according to different criteria, and remove multiple games again. This makes it a great way to manually but quickly narrow down a large dataset according to multiple conditions.