# User Documentation

Samuel González Martín

March 2025

# 1 Software and Hardware Requirements to Run the Project

To run this software, you must first have R installed along with any IDE that allows handling R files. In this case, RStudio has been used, but any other IDE is viable (Visual Studio Code, TinnR, Eclipse, Rbase, etc.). The project has been developed on Windows, so the installation instructions are specifically for this operating system.

Separately, we must also install the MongoDB driver, MongoDB Compass, so the system can create patient records locally. How to install all these programs is detailed in the next section.

# 2 Installation / Setup

## 2.1 Installing R

To install R on Windows, you can download it from the CRAN website [CRAN, 2024]. Once on the website, click on the "*Download*" section, as shown in Figure 1. This will download an executable (.exe) file that will be found in the Downloads folder.

When executing it, you must select the installation language 2. Once selected, the license agreements will appear 3, and you will be given the option to choose the installation directory 4; in this case, the default folder will be used. Then, you will select the components to install and specify if you want to use configuration options, as shown in Figures 5 and 6. Finally, before starting the installation, you will choose a folder to store the program's shortcuts 7 and define additional tasks to be performed during installation 8.
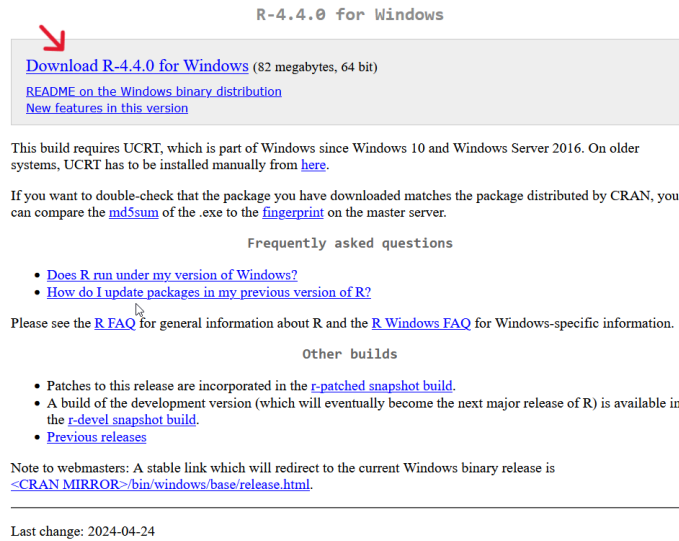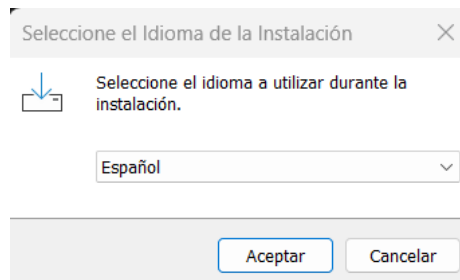
Figure 1: Download R
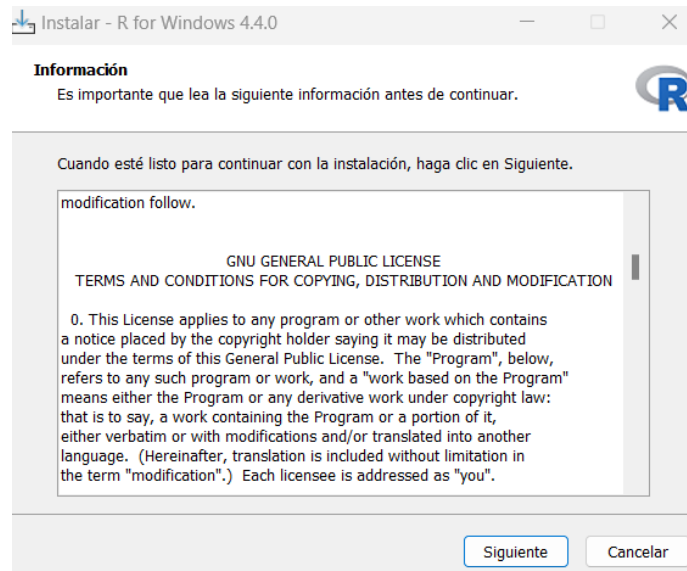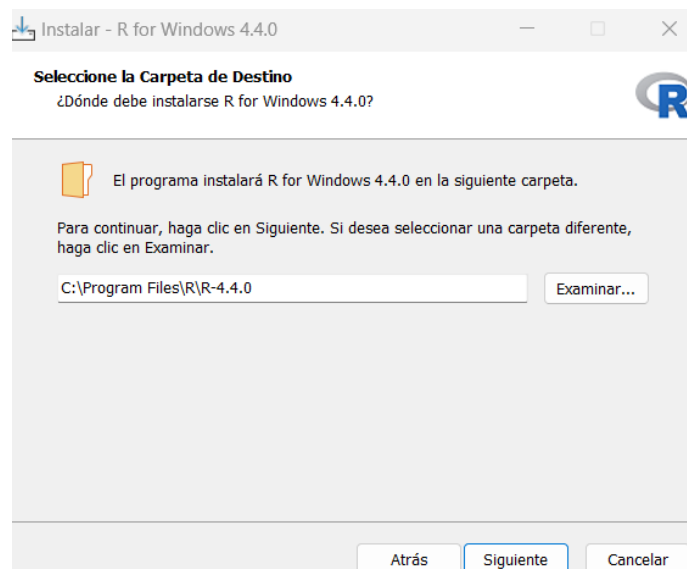


Figure 2: Language Selection

Figure 3: Licenses



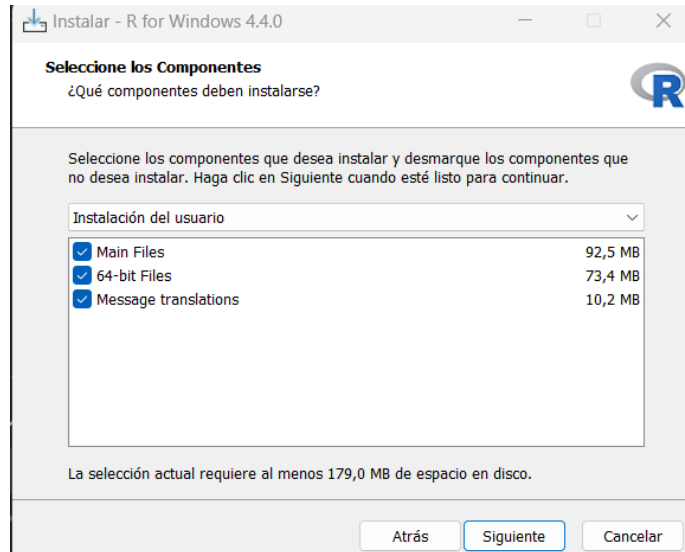Figure 4: Destination Folder Selection
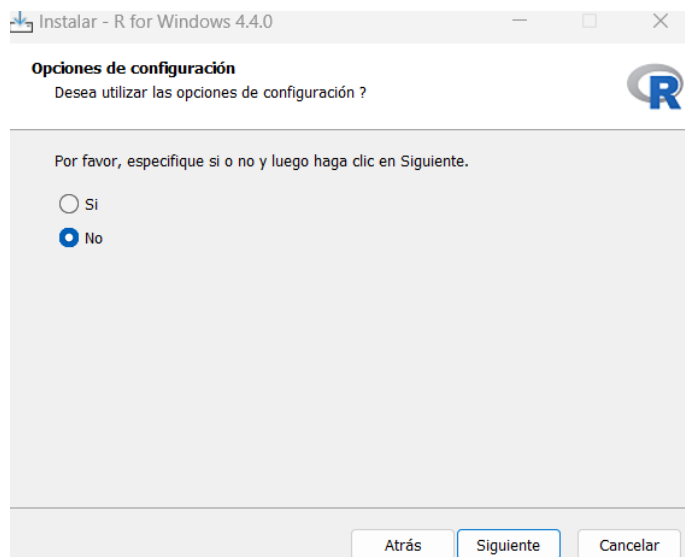
Figure 5: Component Selection
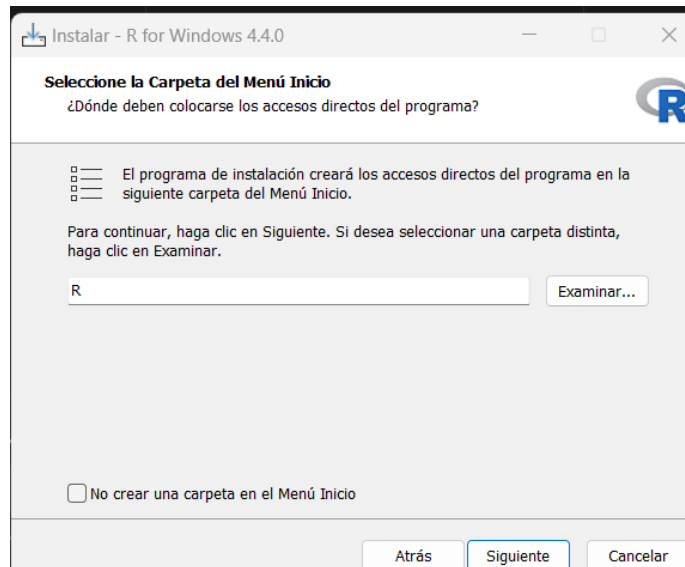


Figure 6: Configuration Options
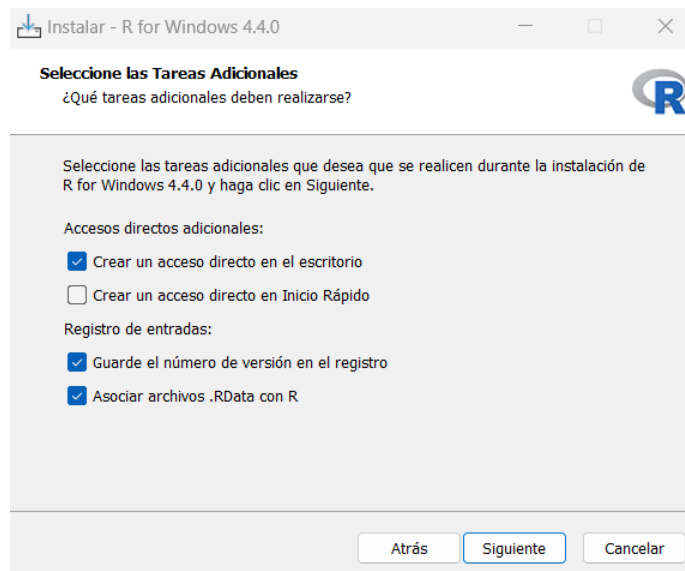
Figure 7: Shortcut Folder Selection



Figure 8: Additional Tasks

## 2.2   Installing RStudio

The installation of RStudio is similar to that of R. To download the installation file 9, go to the Posit website [Posit, 2024].

When executing the file, you will be asked to indicate the destination directory 4 and the folder for RStudio shortcuts 11. Once this information is completed, the installation will begin.



Figure 9: Download Installation File

Figure 10: Destination Folder Selection



Figure 11: Shortcut Folder Selections

## 2.3   Installing MongoDB Compass

On the official MongoDB website [MongoDB, 2024], you can find the MongoDB Compass installer 12. In this case, no configuration is required; the application will install automatically upon executing the downloaded file.



Figure 12: Download MongoDB Compass Installer

Finally, to launch the Shiny application, you only need to clone the GitHub repository, open the PDF-Scrapping.R file, and execute it.

The instructions to clone a repository [GitHub, 2024] are:

- Go to the GitHub repository page.

- Click the green button labeled "¡¿ Code" 13.

- Select the cloning method (HTTPS, SSH, or GitHub CLI).

- Copy the displayed URL.

- Open Git Bash and navigate to the directory where you want to clone the repository.

- Type "git clone" followed by the copied URL.

Figure 13: Repository Cloning

# 3   Manuals and/or Practical Demonstrations

Following the use case diagram in Annex F, demonstrations for each use will be shown.

## 3.1   Processing PDFs

Figure 14 shows how data has been extracted from three uploaded files.



Figure 14: Demonstration of PDF Processing

## 3.2   Querying Searches in Clinical Databases

The column Patogenicidad.Buscada in the table of Figure 15 corresponds to the searches for mutations in the ClinVar database.



Figure 15: Demonstration of querying clinical databases

## 3.3   Modifying Tables

Double-clicking on any cell in the tables will open an editor where we can modify the value of that cell.

Figure 16: Demonstration of table modification

### 3.3.1 Sorting Values

In this example 17, the column Porcentaje_tumoral has been sorted in descending order.



Figure 17: Demonstration of sorting values

### 3.3.2 Filtering Rows

In Figure 18, only records corresponding to male patients have been filtered.



Figure 18: Demonstration of filtering rows

## 3.4 Exporting Data

In this example, the data has been extracted into a CSV document 19, and the visualization of the downloaded file is shown in 20.

Figure 19: Demonstration of exporting data



Figure 20: Demonstration of exporting data

## 3.5 Creating Patient History

In Figure 21, the records of the analyzed documents are shown in the MongoDB Compass application.



Figure 21: Demonstration of creating patient history

# References

[CRAN, 2024] CRAN (2024). Download r-4.4.0 for windows. the r-project for statistical computing. https://cran.r-project.org/bin/windows/base/. Acceso realizado el 5 de junio de 2024.

[GitHub, 2024] GitHub (2024). Clonar un repositorio - documentación de github. https://docs.github.com/es/repositories/creating-and-managing-repositories/cloning-a-repository. Acceso realizado el 11 de junio de 2024.

[MongoDB, 2024] MongoDB (2024). Download mongodb atlas cli — mongodb. https://www.mongodb.com/try/download/atlascli. Acceso realizado el 5 de junio de 2024.

[Posit, 2024] Posit (2024). Rstudio desktop - posit. https://posit.co/download/rstudio-desktop/. Acceso realizado el 5 de junio de 2024.