# Exploratory Data Analysis in Jupyter Notebook inside ULEAD

## ACCESSING THE TOOL AND UNDERSTANDING THE FUNCTIONS AND FEATURES

CENAN PIRANI & JAVIER SANZ

Document version 1.0

# Exploratory Data Analysis Jupyter Notebook

This guide describes how to use a tool for the initial exploration of Electronic Health Record (EHR) data that was generated by the CTSI and delivered to investigators for research purposes. The data for your project has been exported as a set of text files (using the "CSV" format). If you are using the ULEAD enclave, these have been placed in the 'Data' folder for this project.

This toolkit has then been built to contain commands that will describe each of the specific files in your dataset . The software library is written in Python using the Pandas data analytics library on the Jupyter Notebook platform and if you are familiar with Pandas you can proceed to conduct additional analyses in the same notebook.
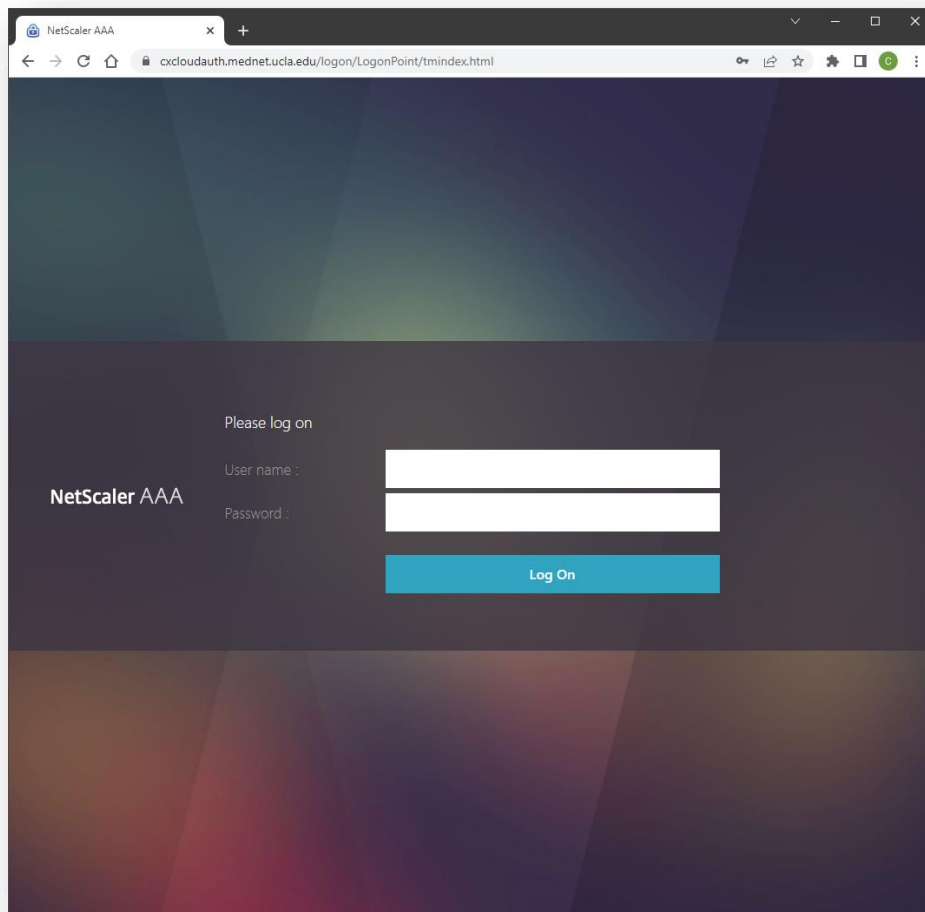
The tool is intended to help Investigators gain a deeper understanding of the data that they are receiving with tools that make it easy to complete basic and some intermediary data analysis tasks. It also aims to educate and encourage Investigators to expand their knowledge of analytics using Python and Pandas.

## THE EDA NOTEBOOK FUNCTIONS

The notebook's main library, 'ehr_dp_lib.py', is made up of a set of functions developed by our team that targets the main descriptive statistics that investigators typically need to begin understanding their data. The code blocks in the notebook are pre-generated to the match the different files and variables contained in your specific data set. The functions are as follows but rather than reviewing these in detail here you can proceed to try examples using your own data:

- **describe_tables():** Returns a dataframe listing all the files in the 'Data' folder including row and column counts and descriptions
- **missingness( dataframe name ):** Returns a dataframe of the number of null values per column.
- **catbar( dataframe name, column name, graph=(True or False))**: [Generated on categorical data type only] Returns a dataframe of counts of all the groups of categories in the specific column in the dataframe. When graph argument set to True returns a bar graph.
- **numstats( dataframe name, column name ):** [Generated on number data type only] Returns a dataframe of descriptive statistics (ie. mean, max, min, median, quartiles) for the column data.
- **dateline( dataframe name, column name ):** [Generated on date data type only] Returns a line graph of the freuency of specific dates along an x-axis of time.

- **flow_stats( flowsheet dataframe ):** [Generated only if Flowsheet_Vitals.csv table in Data folder] Returns a dataframe of descriptive statistics for common vitals sign types (ie. Height, Weight, Temperature, Sp02, Pulse, BMI, Respirations).
- **lab_stats( lab dataframe, top=(10 or greater) ):** [Generated only if Labs.csv table in Data folder] Returns a dataframe of descriptive statistics for top lab procedures in dataset. The top argument can be adjusted to capture more lab procedures.
- **text_search( dataframe name, column name, text to search, ignore case=(True by default can also be set to False) ):** Returns a dataframe based on a free text search of a specific column in an existing dataframe.



## ACCESSING ULEAD WITH ANACONDA AND JUPYTER NOTEBOOK

### STEP 1: Launch ULEAD Citrix Environment

Navigate to: https://uclaohia.cloud.com/Citrix/StoreWeb/

The browser will re-rout and will present a login, enter your AD username / password. Logging in will also prompt the 2FA with Duo.

## STEP 2: Open Anaconda Prompt

Once logged in, you will see the ULEAD Apps screen. If not already listed under "Recent" or "Favorites," click on the link: 'View all applications.'



Then click on the icon 'Anaconda Prompt.'

(you might want to click on the ⭐ next to the icon and mark it as 'favorite' for future sessions).

## STEP 3: Navigate to root folder

Once finished loading a prompt will appear in the window in a virtual screen.

In the prompt enter the following commands:

- F: <enter>
- cd Inbound <enter>
- cd CTSI <enter>

At this point the prompt should look like this:

Now enter the command "dir/w" to list all of the folders in the directory, find your folder and enter: "cd <folder name> <enter>."

## STEP 5: Launch Jupyter Notebook application

Now type the command "jupyter-notebook" and press <enter>. After a short time it will open a Chrome browser with the Jupyter Notebook explorer in your own ULEAD folder.



## STEP 6: Open EDA document

At this point you can click on the file: 'Data_Profiler.ipynb' to launch the notebook: