# seenopsis: a tool for first exploring and visualization of available variables in a dataset

**The Why**

As an epidemiologist, responsible for the design of medical research, once I have a dataset I need to "feel" the variables.

Though essential, this task is sometimes repetitive, time consuming and mainly boring. What if there was a tool that centralizes the main important features of all the variables in the dataset, helping you explore the variable in a structured visualized approach?
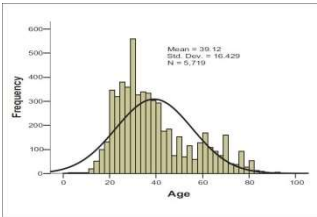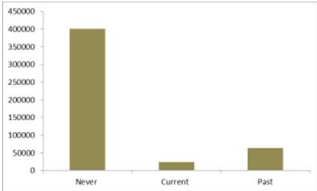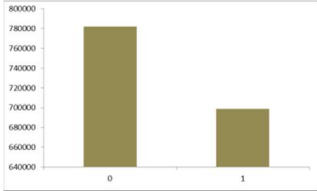
**The What**

**seenopsis** is designed to help the everyday work of a data scientist, to first explore the available variables in a dataset.

**The How**

The only required argument in **seenopsis** is the name of the dataset. Other arguments are optional.

**The Vision**

This is how I imagine the output of s**eenopsis:**

| Variable name | Type of variable | Distribution | Missing | Basic stat | Outliers |
|---|---|---|---|---|---|
| **age** | Continuous |  | Missing: 31,790 Total population: 3,695,863 Percent of missing: 0.86% | Min: 0.01 Max: 102 Mean: 40.3 Median: 38.5 SD: 1.1 | >1.5 sd: 80 >2 sd: 50 >3 sd: 3 Optional [values, direction] |
| **Smoking_ status** | Categorical (3 categories) |  | Missing:1,145,717 Total population: 3,695,863 Percent of missing: 31.0% | - | - |
| **Stroke_ comorbid** | Binary |  | Missing: 609 Total population: 3,695,863 Percent of missing: 0.02% | Min: 0 Max: 999 | - |