Generating Automated Dataset Summaries

Agustin Calatroni and Herman Mitchell, Rho, Inc.

Overview

Every time a dataset is created, either for data management purposes or for statisti- We have a SAS dataset and we need to create a statistical summary for all the varical analyses, it is imperative that each variable be reviewed carefully. A summary re- ables. The SAS macro insert shows all the current options in the macro. We specify ess appears, there is no straight forward procedure in SAS to produce such a report. grams to generate an automated summary report of dataset variables.

Motivation

- Show all the data; display many numbers in small spaces
- Reveal the data at several levels of detail, from broad overview to detailed structure
- Serve as a description, tabulation and exploration of the data

Objective

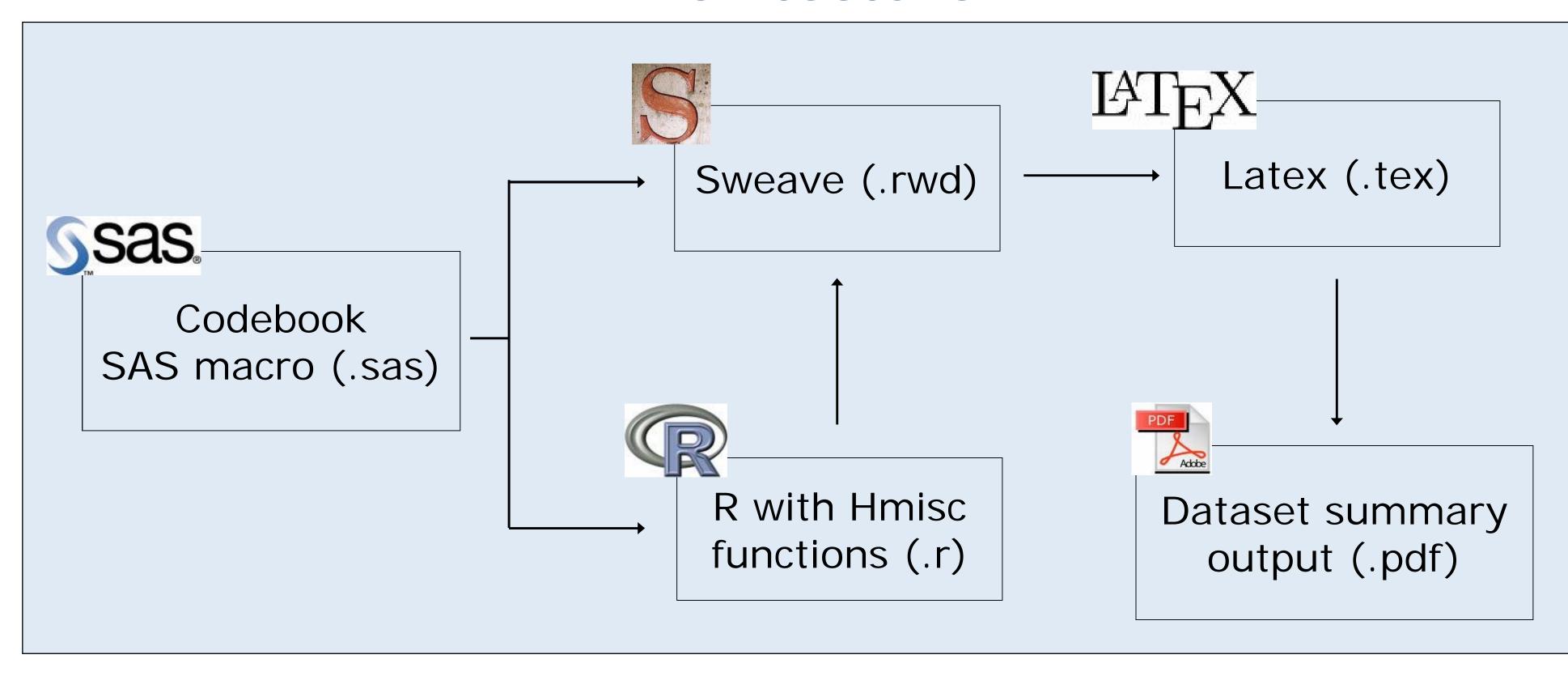
Create a concise statistical description of each variable in a dataset with the following characteristics:

- Output should be dictated by the variable type
- Output should be able to combine text and graphics (show all the data, including extreme values - hide nothing!)
- A statistical summary should be able to run in SAS using SAS datasets (SAS is the lingua franca for both pharma and CROs)
- A statistical summary should be automated so that it requires minimum input

Important Considerations

- Must include variable labels and formats
- Allow for special numeric missing values
- Display dates and times in a legible way
- Ability to subset variables and/or observations

Architecture



How the SAS Macro Works

The SAS macro works by interacting with a variety of programs to obtain the final dataset summary output. The macro creates two files:

- A Sweave file with:
- Chunk to import the SAS dataset
- Chunk to generate the latex descriptive statistics
- An R file with code to submit the code and generate the output

Guided Tour

port of the dataset should succinctly display critical information to enhance its ease of the format library, pdf location and pdf name for the final output. The macro creates examination and allow dissemination among interested parties. As basic as this proc-two files. The Sweave file contains the R code in chunks intertwined in Latex. You can distinguish the R code because it is enclosed between <<>>= and @. The second file It is of interest to develop a SAS macro that combines the capabilities of several pro- is a set of R commands that will submit the Sweave file and create the pdf from the latex file. The SAS macro executes the set of R commands in batch mode using windows commands from within SAS (X command). From the analyst's standpoint none of these steps are visible, only the SAS macro call and the final pdf output.

> A strength of the current implementation of this dataset summary generator is that the data displays depends upon the variable type. The final output will contain variable names, labels, formats, observations (n), frequency of missing values (nmiss), and number of unique observations for all variables. Special missing codes (._,. ,.a,.b,...,z) will be printed when appropriate. Then depending on the variable type: sum, mean, sd, quantile, frequency and the five lowest and highest values will be presented. When the variable is numeric and there are more than twenty unique values, output will also include a spike histogram showing all the data.

Installation Steps

Although the SAS macro is built to interact with other programming languages automatically, the installation requires some care:

- . Obtain the SAS macro codebook.sas
- 2. Download and install R under c:\
- Add Hmisc packages by Frank Harrell Jr.
- Obtain extra R functions from the author
- 3. Download and install Latex or Mitex if working under Windows (with add-on usepackages setspace, relsize and fancyhdr)
- . Call SAS macro codebook.sas

Enhancement

I'm currently working on two enhancements:

- The ability to run the additional applications from a server to alleviate the installation of additional programs.
- An option that will adapt the output to the case of repeated data

Pros

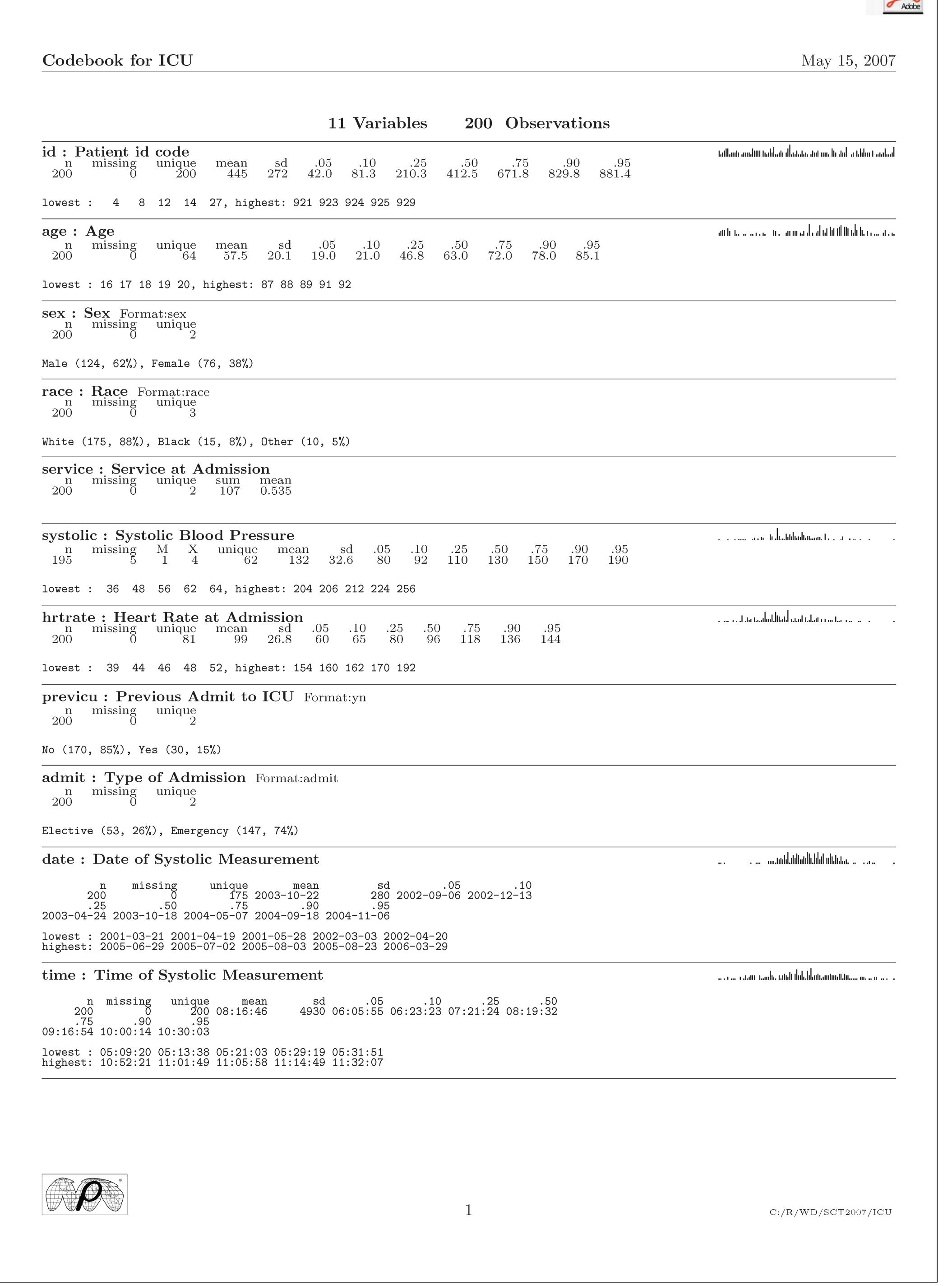
- Highly adaptable program, possibility to include a wide range of summary statistics, plots, etc.
- Attractive displays for data cleaning and distribution analysis
- Extremely user friendly macro once installed

- Installation of additional programs requires extra care
- SAS is used as an interface to other programming languages
- Further customization requires complex programming in R and Latex

Conclusion

Although the current installation may be slightly intricate, the resulting output greatly simplifies data cleaning and analysis.

Final dataset summary output



Terminology and Sample Code

Sweave is a flexible framework for mixing text and R code for automatic document generation. A single source file contains both documentation text and R code, which are then intertwined into a final document containing the combined text and R code and/or the output of the code (text, graphs). What is R?

R is a language and environment for statistical computing and graphics. R is an implementation of the S programming language. R provides a wide variety of statistical and graphical techniques. The capabilities of R are extended through user-submitted packages, which allow specialized tasks. The Hmisc package developed by Frank Harrell Jr. is one such addition. It offers specialized functions to convert SAS dataset to R objects, create concise statistical summaries and conversion to Latex.

What is Latex?

Latex is a document preparation system for high-quality typesetting. It is often used for medium to large technical or scientific documents because of its excellence at managing equations, figures, tables and indices. As with R, Latex is highly extensible through add-ons.

