

Application and available datasets for the analysis of distributional SD

A. Irpino, R. Verde
ESTP Cologne 14-16 May 2024

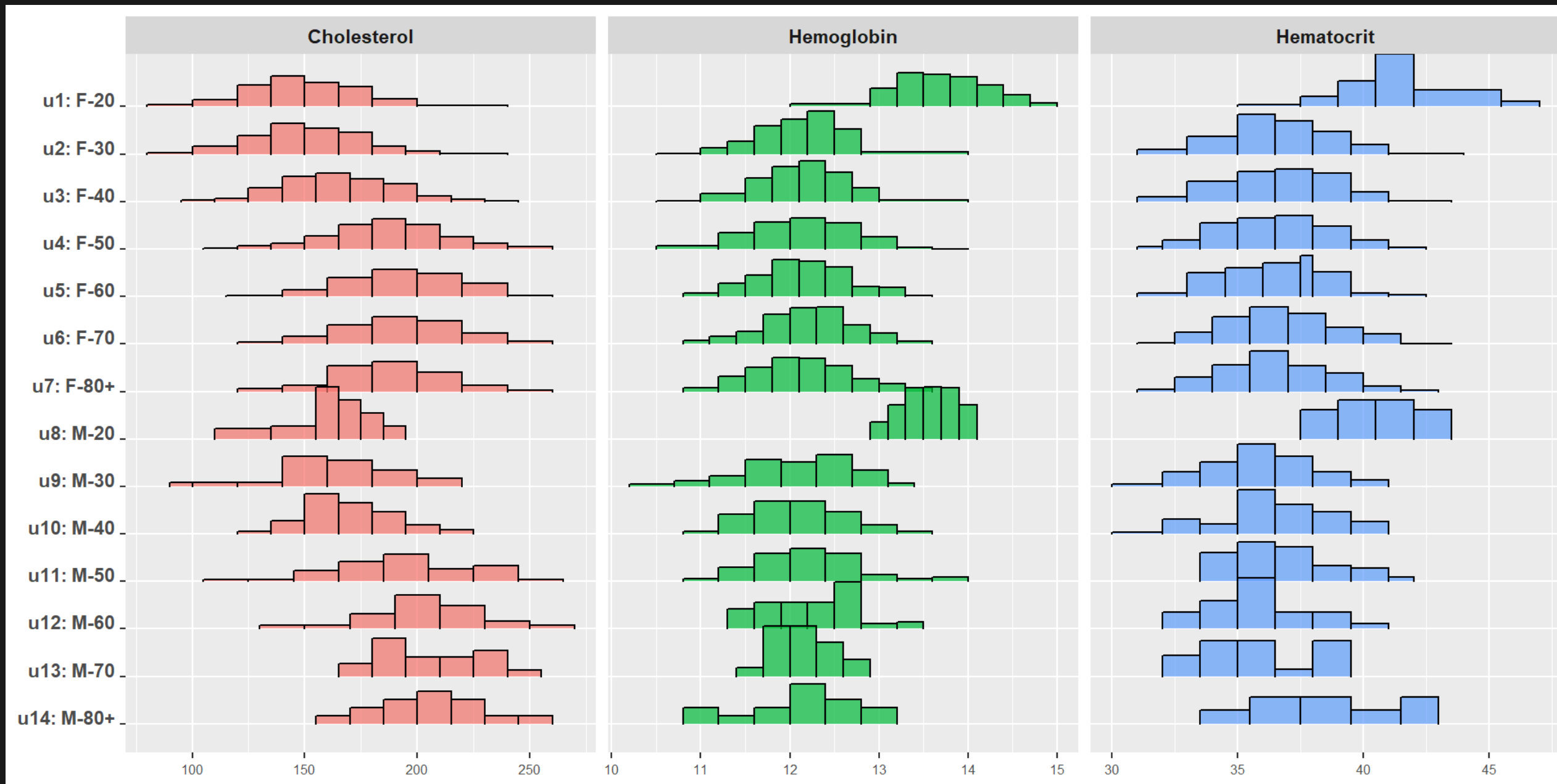
Datasets available for exercising

Datasets contained in the **HistDAWass** package

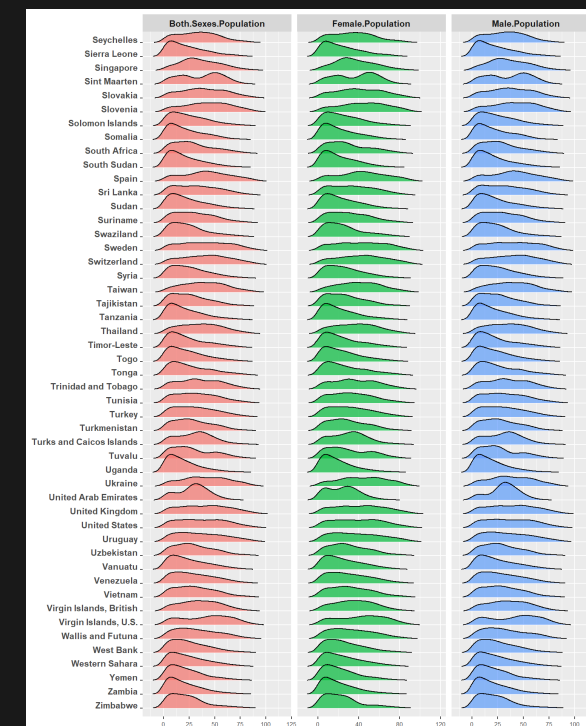
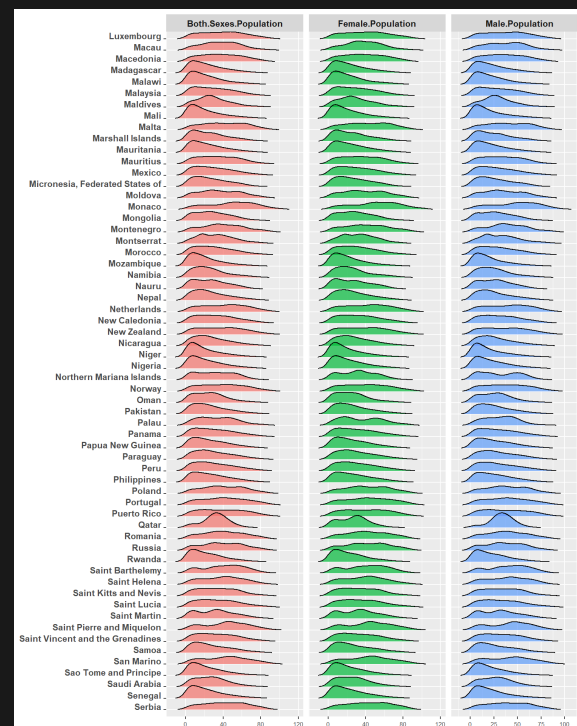
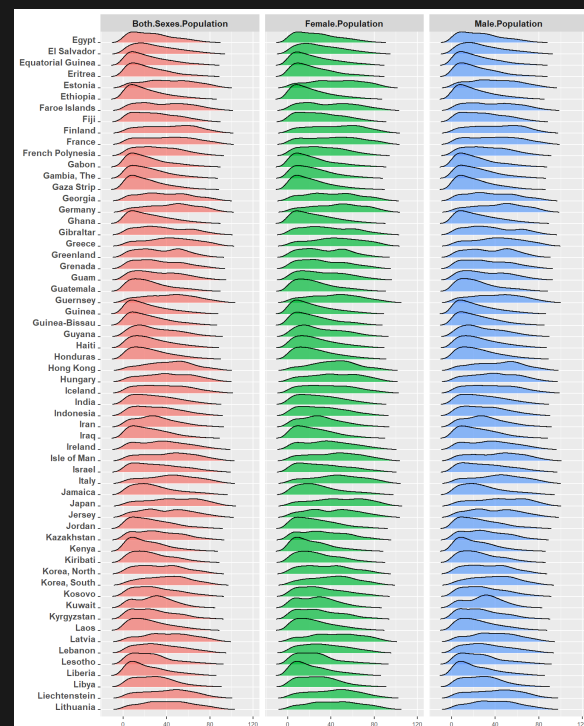
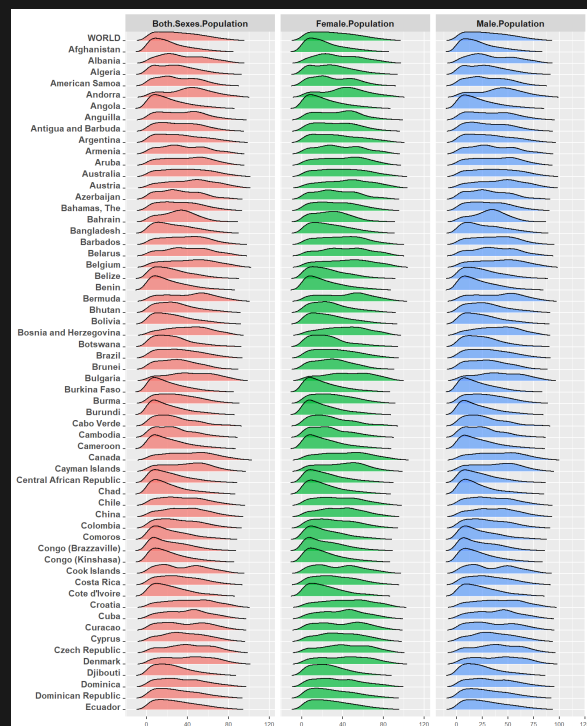
The **HistDAWass** package in R contains some distributional datasets generated from some books and from raw data from different contexts of application. Let's see some of them:

- **BLOOD** dataset, is the description of 14 typologies of patients in a hospital, the variables are histograms of the levels of Hematocrit, Hemoglobin and Cholesterol
- **Age_Pyramids_2014** Age pyramids of all the countries of the World in 2014
- **China_Month** Distributions of climatic variables for each month of 60 stations
- **China_Seas** Distributions of climatic variables for each season of 60 stations China
- **OzoneFull** Full Ozone dataset for Histogram data analysis

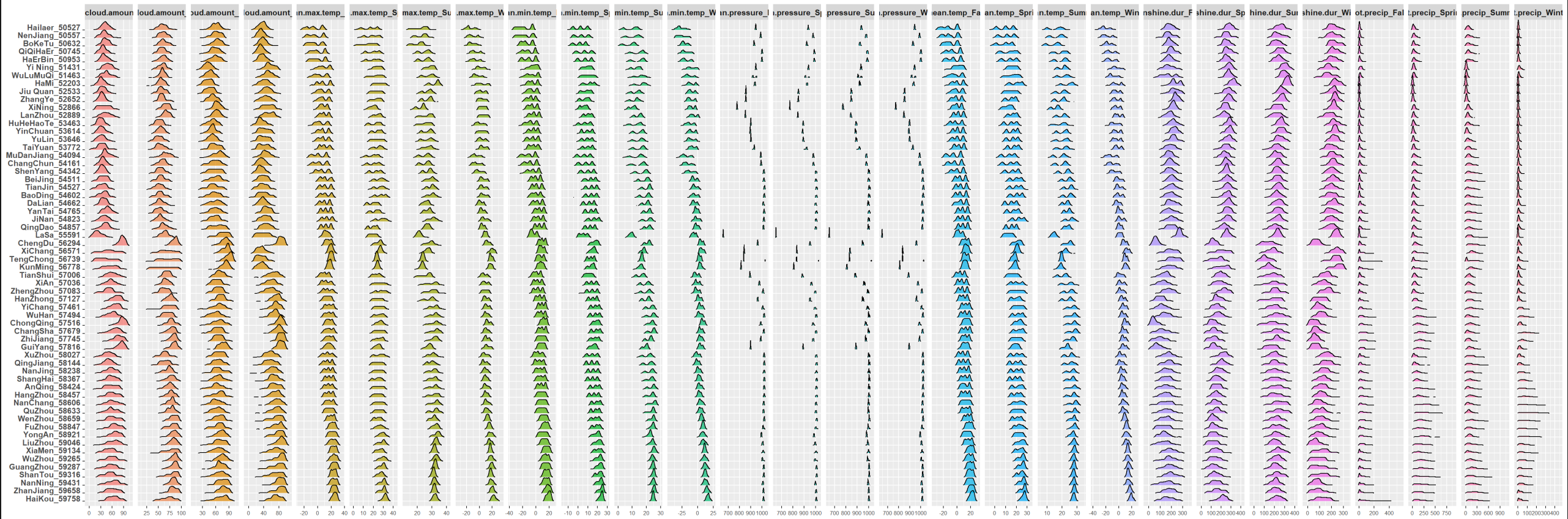
The BLOOD dataset

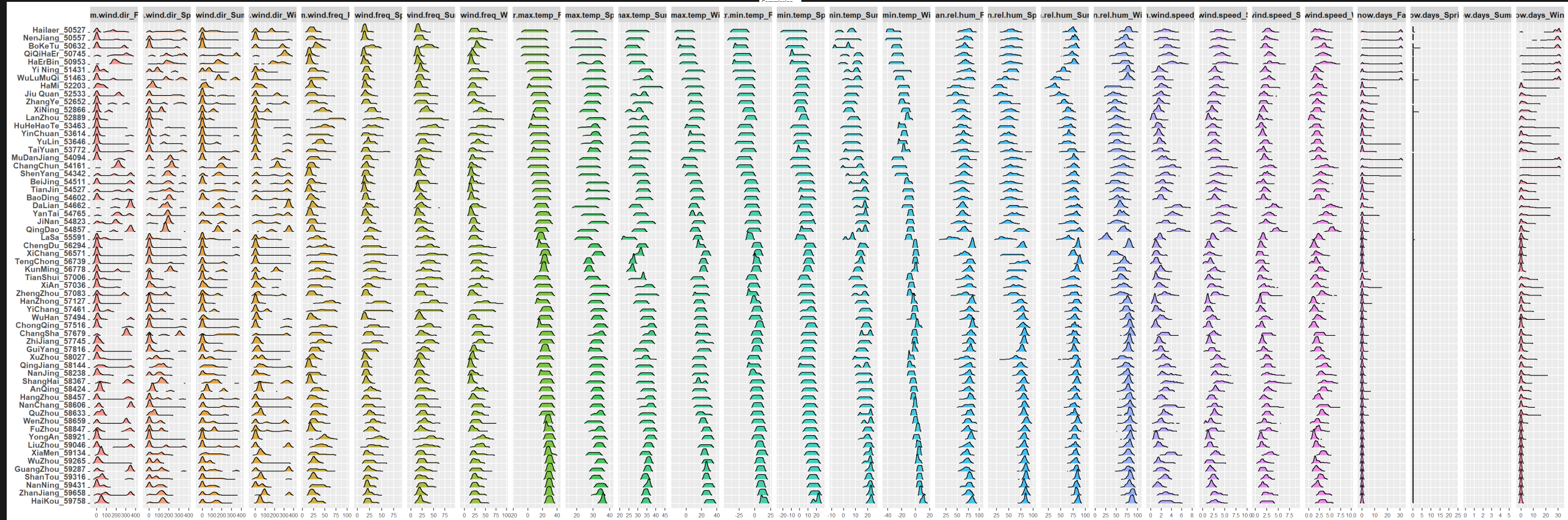


The AGE PYRAMIDS dataset (229 countries)

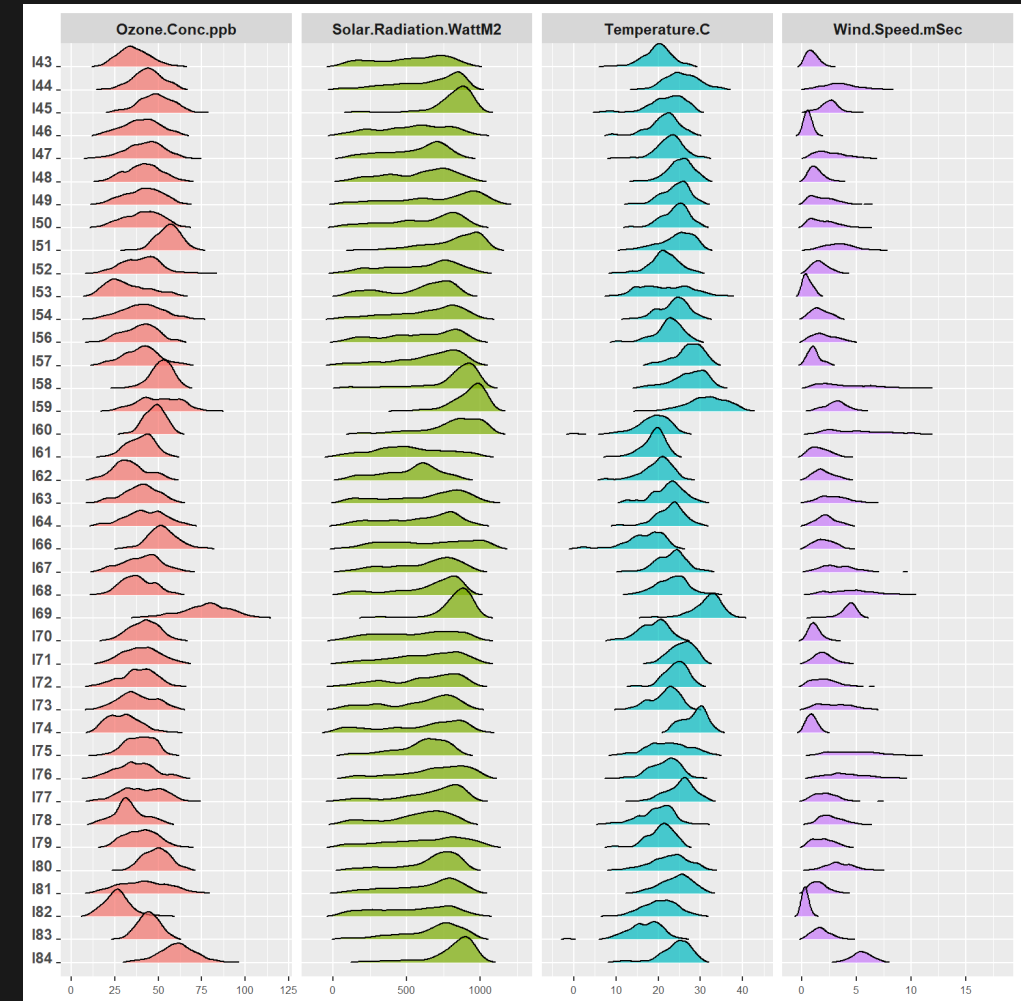
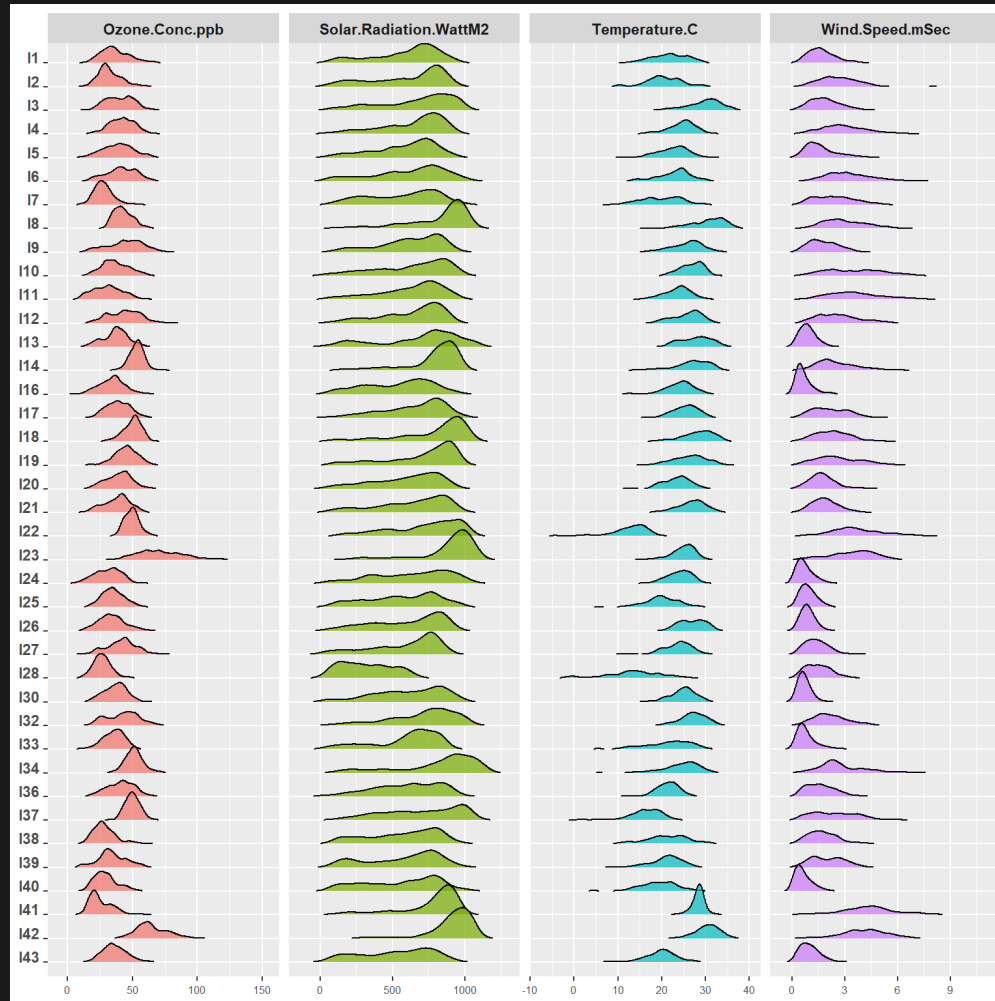


The CHINA season dataset 60 obs, 56 vars



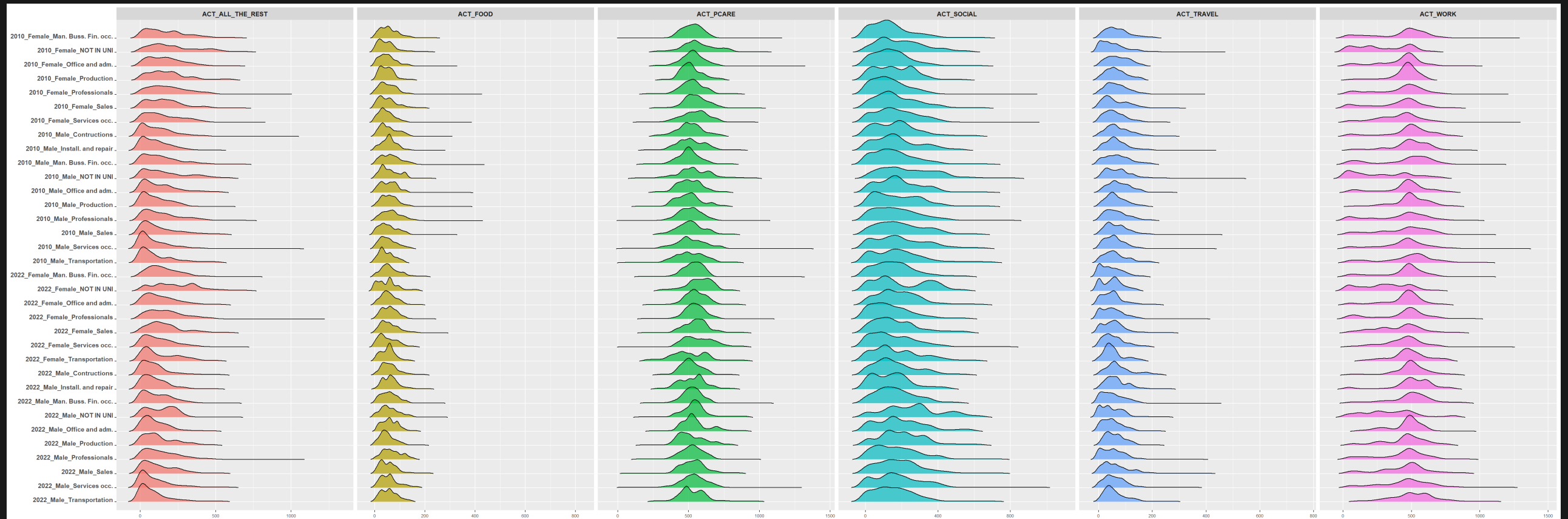


The OZONE dataset 78 obs, 4 vars



Other distributional datasets

- **HMAT** Time use distributional data created from from an IPUMS dataset available from https://github.com/Airpino/Symbolic_data_analysis_softw/raw/main/Histo_data_IPUMS.RData



Some exploratory tools

Some links showing some exploration tools for distributional data

- An application on climatic data: https://airpino.github.io/ICDS_23_presentation/#/title-slide
- An application on EUSILC 2013 data: https://airpino.github.io/DSSR2024_pres/#/title-slide

An analysis task

1. Choose a dataset
2. Explore the data
 - compute some basic statistics from data
3. Fix some objectives
4. Analyze the data (or a selection of the data) using at least one of the following techniques
 - Principal components analysis
 - Regression analysis
 - Clustering analysis
5. Present the results
6. Let's discuss together