

NOTE 1. INTRODUCTION

INTRODUCTION TO STATISTICAL PROGRAMMING

Chanmin Kim

Department of Statistics
Sungkyunkwan University

2022 Spring

COMPUTING IN STATISTICS

- Statistics = Mathematics + Computer Science.
- Size of data $\uparrow \Rightarrow$ Data analysis.
- Complex methods \Rightarrow Numerical/Sampling methods.
- Big data \Rightarrow Parallel/Cloud computing.
- Computationally intensive methods.

STATISTICAL SOFTWARES

- SAS: Commercial and analytic software (Business solution).
- SPSS: Commercial software for social science.
- R: Free software for statistical computing & graphics.
- Minitab, S-plus, Stata, JMP, Octave, etc.

IF YOU GET AN INDUSTRY JOB RELATED TO STATISTICS,

- Data handling:
 - ▶ Extracting data from database.
 - ▶ Missing data.
 - ▶ Changing variables.
 - ▶ Creating new variables.
 - ▶ Splitting / Merging datasets, etc.
- Data analysis:
 - ▶ Descriptive statistic.
 - ▶ Applying models to data.
 - ▶ Visualization.
- Report: Summarizing results.

IF YOU GET AN ACADEMIC JOB,

- Complex / unstructured data: Correlated data, functions, images, text.
- Complex models: No analytical solution \Rightarrow Computational methods.
- Computational statistics: Development of Fast/efficient algorithms.
- Simulation.

HISTORY OF R

- S language:
 - ▶ developed as an internal statistical analysis environment by AT&T lab in 1976.
 - ▶ S-plus (Commercial statistical software).
 - ▶ Data analysis.
 - ▶ Interactive environment.
 - ▶ For user & developer.
- R language:
 - ▶ Created in 1991.
 - ▶ Similar syntax to S.
 - ▶ Free software.
 - ▶ The R Core Group formed in 1997 controls the source code for R.

FEATURES OF R

- Available in various operating systems (e.g., Windows, Mac, Linux, etc.) (\therefore open source nature).
- Sharing with many popular open source projects \Rightarrow active development \Rightarrow frequent releases.
- Statistical analysis purpose + General purpose programming.
- Interactive environment.
- Object-oriented & functional programming language.
- Active & vibrant user community \Rightarrow Development of platforms or packages.

OBJECT-ORIENTED & FUNCTIONAL PROGRAMMING

- Object-oriented programming:
 - ▶ You can pick and choose parts of an object.
 - ▶ Polymorphic \Rightarrow Generic function (e.g., `plot()` can be used for both data objects and regression objects).
- Functional programming:
 - ▶ Avoidance of explicit iteration.
 - ▶ More compact code.
 - ▶ Potentially much fast execution speed.
 - ▶ Less debugging time.
 - ▶ Easier transition to parallel programming.

R SYSTEM

- Base R system: Download from CRAN (Comprehensive R Archive Network).
- Packages:
 - ▶ Over 10,000 packages on CRAN.
 - ▶ Many packages from Bioconductor project.
 - ▶ Personal packages.

LIMITATION OF R

- Memory problem: R objects must be generally stored in physical memory \Rightarrow More memory hog than other statistical software.
- Relatively slow speed, especially for loop statement.
- Functionality depends on customer demand and voluntary user contribution.