# Introduction to R

**Setia Pramana**
**Politeknik Statistika STIS**

fppt.com

# Available Statistical Packages

fppt.com

# Statistical Software Used



Setia Pramana

3

# Python or R?

| | Python<br>King of Data Science Programming Languages | R<br>Golden Child of Data Science |
|---|---|---|
| **PURPOSE OF EXISTENCE** | Python is a general purpose multi-paradigm programming language for data science that has gained wide popularity-because of its syntax simplicity and operability on different eco-systems. | R is an open source programming language and environment for statistical computing and graphics available on Linux, Windows and Mac. |
| **USABILITY** | Python language makes it easy for programmers to write maintainable, large scale robust code. | R language has array-oriented syntax making it easier for programmers to translate math to code, in particular for professionals with minimal programming background. |
| **FEATURES** | • OPEN SOURCE<br>• BROADNESS<br>• EFFICIENT<br>• CAN BE EASILY MASTERED UNDER EXPERT GUIDANCE-READ IT, USE IT WITH EASE<br>• EXTENSIBLE | • OPEN SOURCE<br>• ALL-IN-ONE PACKAGE OF A STATISTICAL<br>• ANALYSIS TOOLKIT<br>• EXCELLENT CHARTING BENEFITS ROBUST AND VIBRANT ONLINE COMMUNITY<br>• POWERFUL PACKAGE ECOSYSTEM |

| | | |
|---|---|---|
| **SALARY** | 2014 DICE TECH SALARY SURVEY AVERAGE SALARY FOR PYTHON PROGRAMMERS IS<br>**$94,139** | 2014 DICE TECH SALARY SURVE AVERAGE SALARY FOR R PROGRAMMERS IS<br>**$115,531** |
| **LIBRARIES & PACKAGES** | • NUMPY/SCIPY<br>• PANDAS<br>• SCIKIT-LEARN<br>• STATSMODELS<br>• MATPLOTLIB | • CARET<br>• GGVIS,GGPLOT2<br>• STRINGR<br>• ZOO<br>• PLYR,DPLYR |
| **APPLICATIONS** | • WALT DISNEY USES PYTHON LANGUAGE TO ENHANCE THE SUPREMACY OF THEIR CREATIVE PROCESSES.<br><br>• DROPBOX IS COMPLETELY WRITTEN IN PYTHON LANGUAGE WHICH NOW HAS CLOSE TO 150 MILLION REGISTERED USERS.<br><br>• PYTHON PROGRAMMING IS USED BY MOZILLA FOR EXPLORING THEIR BROAD CODE BASE. MOZILLA RELEASES SEVERAL OPEN SOURCE PACKAGES BUILT USING PYTHON. | • FORD USES OPEN SOURCE TOOLS LIKE R PROGRAMMING AND HADOOP FOR DATA DRIVEN DECISION SUPPORT AND STATISTICAL DATA ANALYSIS.<br><br>• ZILLOW MAKES USE OF R PROGRAMMING TO PROMOTE THE HOUSING PRICES.<br><br>• INSURANCE GIANT LLOYD'S USES R LANGUAGE TO CREATE MOTION CHARTS THAT PROVIDE ANALYSIS REPORTS TO INVESTORS. |

# Python or R?

- [https://www.datacamp.com/community/tutorials/r-or-python-for-data-analysis](https://www.datacamp.com/community/tutorials/r-or-python-for-data-analysis)

# What is R?

- A language and environment for statistical computing and graphics.

- An integrated suite of software facilities for data manipulation, calculation and graphical display.

- First appeared in 1996 by Prof. Ross Ihaka and Robert Gentleman of the University of Auckland, NZ.

- GNU software -> Free. Similar like S language.

- Open source, maintained and developed by a community of developers.

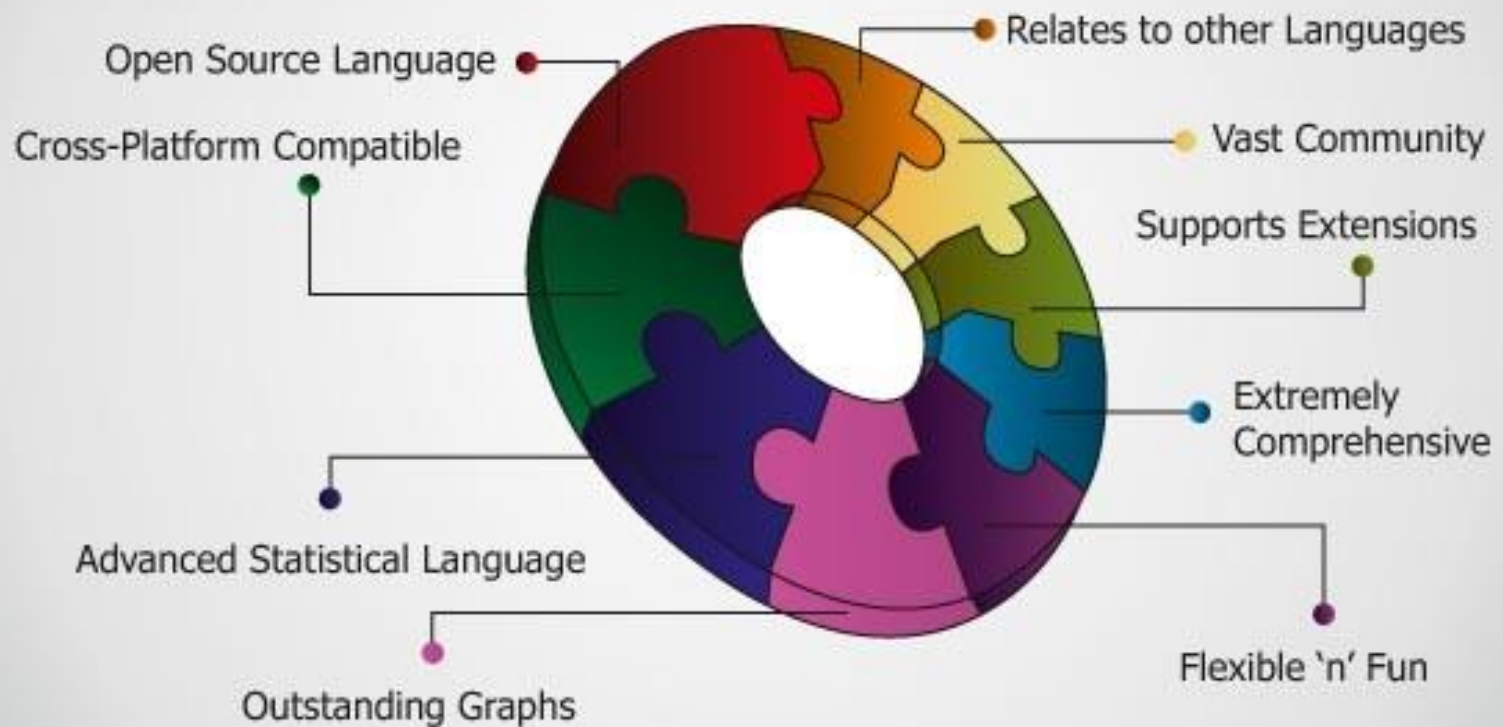- Works in Windows, Unix, MacOs

fppt.com

# R Includes

- Effective data handling and storage facility,
- A suite of operators for calculations on arrays, in particular matrices
- A large, coherent, integrated collection of intermediate tools for data analysis,
- Graphical facilities for data analysis and display either on-screen or on hardcopy
- Well-developed, simple and effective programming language which includes conditionals, loops, user-defined recursive functions and input and output facilities.
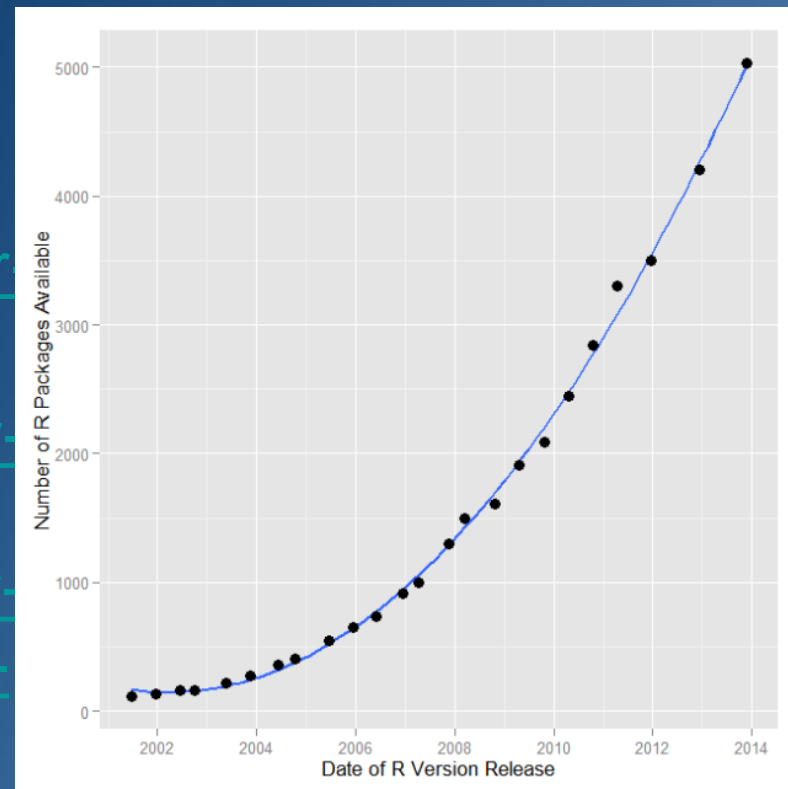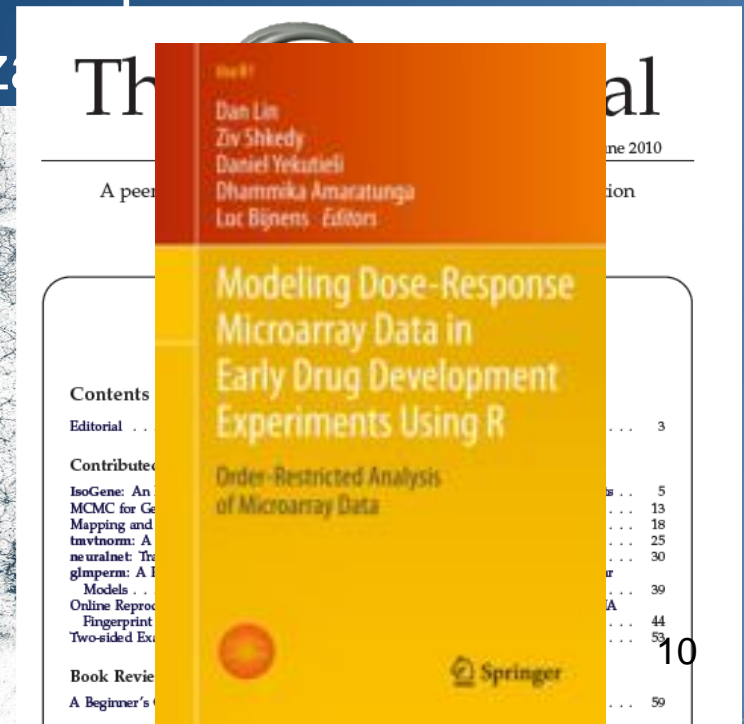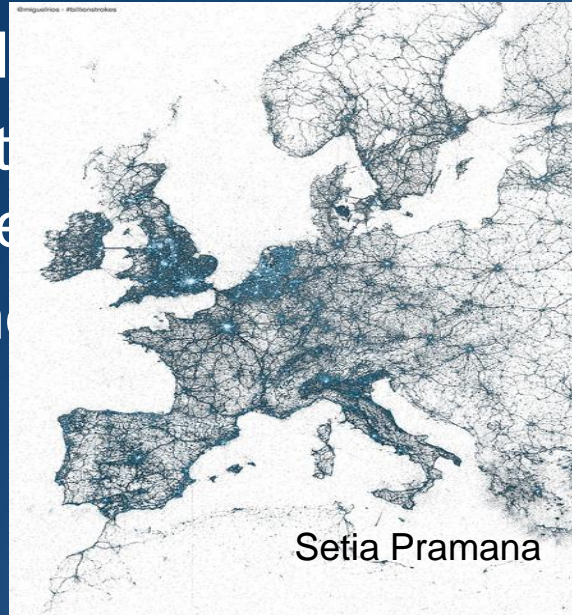
http://www.r-project.org/

fppt.com

# Why R?

- It is not only statistical software but also a language
- 10.000 add-on packages → lots of pre-prepared packages (http://cran.r-project.org/web/packages/)
- With many applications http://cran.r-project.org/web/views/, http://www.revolutionanalytics.com/r-language-features-applications-and-extensions#thirdparty .
- Access to powerful, cutting-edge analytics



Setia Pramana

9

# Why R?

- Flexible (complex or standard statistical practices, bayesian modelling, GIS map building, building interactive web applications, building interactive tests, etc.   )

- We can make our own package and publish it

- Great Graphics and data visualiza

- Can be used

- Well Support
  resources-we

- And many m

Setia Pramana

10

# Why R?

- Can be integrated with other languages (C/C++, Java).

- R can interact with many data sources and other statistical packages (SAS, Stata, SPSS, and Minitab).

- For the high performance computing task → multiple cores, either on a single machine or across a network.
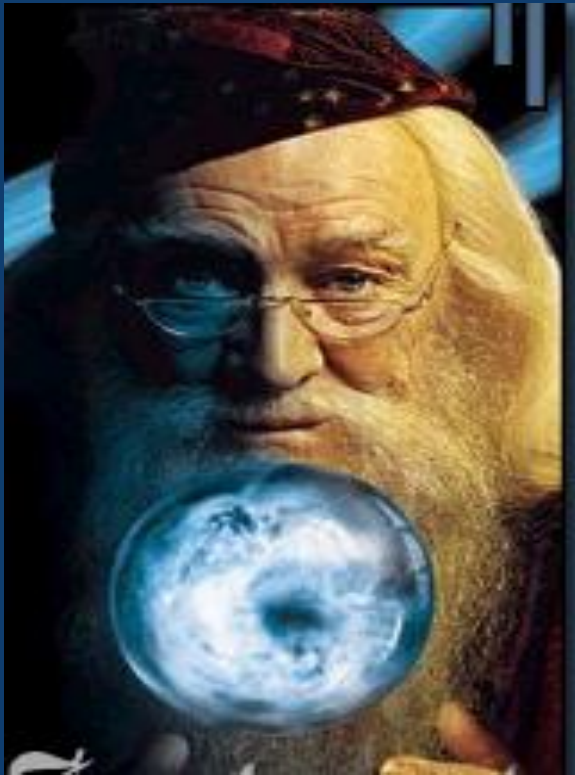
# But…..

- R has no warranty

- Command Line Interface : difficult for some users.

- Users must learn a new way of thinking about data and data analysis sequence
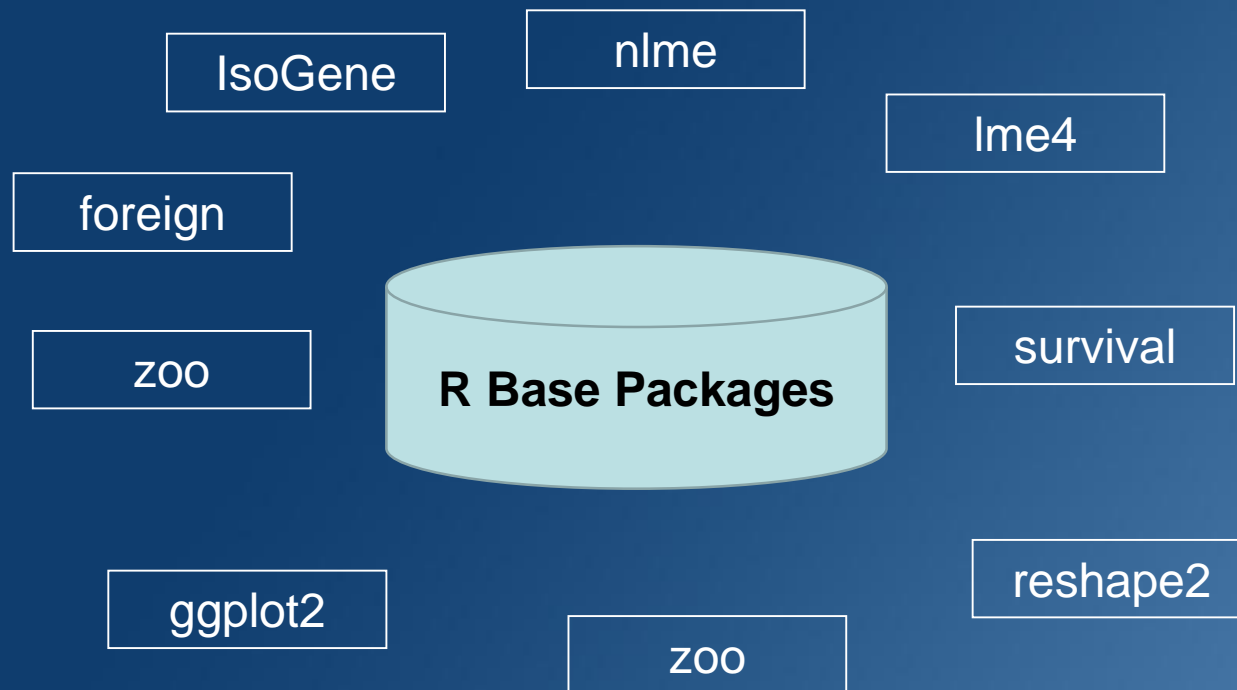
- That's all ….. I guess

# WizaRd

# Learning R

Companies that use R for Analytics

# R Library/packages

IsoGene

nlme

lme4

foreign

**R Base Packages**

zoo

survival

ggplot2

zoo

reshape2

# My R Packages

- IsoGene
- IsoGeneGUI
- nea
- neaGUI
- biclustGUI
- OCRME
- More detail: http://setiopramono.wordpress.com/r-programming/

# Cutting Edge Technologies

fppt.com

# Finance

- https://cran.r-project.org/web/views/Finance.html

- zoo - Provides the most popular format for saving time series objects in R.

- xts - Very flexible tools for manipulating time series data sets.

- quantmod - Tools for downloading financial data, plotting common charts, and doing technical analysis.

- The QuantTools package offers enhanced quantitative trading and modeling tools.

- The Risk package computes 26 financial risk measures for any continuous distribution.

# Data Mining

- Decision trees: rpart, party
- Random forest: randomForest, party
- SVM: e1071, kernlab
- Neural networks: nnet, neuralnet, RSNNS
- Performance evaluation: ROCR
- Data Mining GUI rattle
- etc..
- http://www.rdatamining.com/

# Social Media

- Text mining: tm

- Topic modelling: topicmodels, lda

- Word cloud: wordcloud

- Facebook: RFacebook

- Twitter data access: twitteR

- Social Network: sna, igraph, RSiena

    (http://www.jstatsoft.org/v24/i06/paper)

    http://www.r-bloggers.com/an-example-of-social-network-analysis-with-r-using-package-igraph/

# Parallel Computing

- snow (Simple Network of Workstations) & snowfall for development of parallel R programs.

- multicore parallel processing of R code on machines with multiple cores or CPUs

- More: http://cran.r-project.org/web/views/HighPerformanceComputing.html

# Big Data

- <u>RHadoop</u> - a collection of five R packages that allow users to manage and analyze data with Hadoop, developed by Revolution Analytics

- <u>RHIPE</u> - an R and Hadoop Integrated Programming Environment

- More…...

# R Graphics and Visualization

- R provides wide range graphics and visualizations
- Basic Plots: bar plots, basic 3D plots, heatmap.,etc
- Geographic Maps
- Projection Maps
- Social Network Graphs
- Animated graphics and movies (animation)
- Motion Charts (GoogleViz)
- Interactive Graphics (rggobi)
- Image format: BMP, JPEG, PDF, PNG etc…
- More: https://www.r-graph-gallery.com/

Setia Pramana

24

# R Graphics

# R Graphics



Based on estimates from:
Abel and Sander (2014) *Science* Vol. 343 no. 6178 pp. 1520 – 1522

# R Graphical User Interfaces

- R uses Command line interface  and it is preferred for advanced users → allows direct control, more accurate, flexible and the analysis is reproducible.

- Requires good knowledge of the language → difficult for beginners or less frequent users.

- R provides tools for building GUIs → RGUI

# R GUI Projects

- Integrated development environment (IDE)/Script Editors aimed to provide feature-rich environments to edit R scripts and code: Rstudio ([www.rstudio.com](www.rstudio.com)), and architect ([www.Openanalytics.eu](www.Openanalytics.eu))

- Web based application: the Rweb (Banfield, 1999), R.Net ([www.u.arizona.edu/~ryckman/Net.php](www.u.arizona.edu/~ryckman/Net.php)), or gWidgetsWWW (Verzani, 2012).

# R GUI Projects

- Python: **OpenMeta-Analyst** (Wallace et al, 2012)
- Java: **JGR** (Java GUI for R), **Deducer** (Fellows, 2012), and **Glotaran** (Snellenburg, 2012).
- Php: R-php ([http://dssm.unipa.it/R-php/](http://dssm.unipa.it/R-php/))
- Other extensions connect R to graphical toolboxes for developing menus and dialog boxes: **Tcltk, Gtk.**

Setia Pramana

# R Studio

- Download from Rstudio.com
- Powerfull IDE (Integrated Development Environment) for R.

# RGUI Developed using tcltk

# RGUI: RCommander

- [Rcommander.com](Rcommander.com)
- Helpful for R beginner
- Install inside R

# RGUI: Web Based App

# WebBUGS

- Conducting Bayesian Statistical Analysis Online

- Combines OpenBUGS and R

www.webbugs.psychstat.org

# RGUI: Shiny

- A new package from Rstudio to build interactive web applications with R.

- Really Easy!

- Build useful web applications with only a few lines of code—no JavaScript required.

- Self learning: http://shiny.rstudio.com/

- http://www.showmeshiny.com/

# Our Recent R Packages

| Name | Title | Brief Description | Author | Repository |
|------|-------|-------------------|--------|------------|
| spatialClust | Spatial Clustering | Clustering analysis with pay attention on membership via spatial effects | Imam Habib Pamungkas, Setia Pramana | CRAN |
| advclust | S4 Object Oriented for Advanced Clustering( Fuzzy Clustering and Cluster Ensemble) | Advance on clustering with fuzzy clustering for overlapping cluster and objects on gray area. Cluster Ensemble performs combining several result as one robust and stable result. | Achmad Fauzi Bagus F, Setia Pramana | CRAN |
| RcmdrPlugin.Fuzzy Clust | R commander plugin for fuzzy clustering | Graphical User interface via Rcmdr Plugin for fuzzy clustering analysis | Achmad Fauzi Bagus F, Setia pramana | CRAN |
| MetaheuristicFPA | Metaheuristic with Flower Pollinantion Algorithm | Optimization of function objectives to get global optimum of parameter by using Flower Pollination Algorithm | Amanda Pratama Putra, Margaretha Ari Anggorowati | CRAN |
| Multiplier | Social Accounting Matrix and Finansial Social Accounting Matrix | Graphical User Interface for performing SAM (Social Accounting Matrix) and FSAM (Financial Social Accounting Matrix) | Tiara Ratna Dewi, Aisyah Fitri Yuniarshi | R-Forge |
| RcmdrPlugin.PCAR obust | Robust PCA plugin for Rcmdr | Graphical User Interface for Robust Principal Component Analysis (PCA) with Hubert Algorithm for Dimension Reduction | Monalisa Sipahutar, Setia Pramana | CRAN |

fppt.com

# Our Recent R Packages

# Our Recent R Packages

- Kalingga
- Muria
- C++

**Asgard** Alpha Version

ASGARD is a statistics software used to perform geographically weighted regression (GWR). This software was made in 2016 and currently contains some basic GWR functions like GWR, Geographically Weighted Poisson Regression (GWPR), Geographically Weighted Logistic Regression (GWLR), Geographically Weighted Negative-Binomial Regression (GWNBR) and some Assumption Test related to GWR. In addition, ASGARD is also integrated with the map that make it easier for users to performs analysis.

## MAIN FEATURES

Spreadsheet

Fairly complete functions
· GWR
· GWPR
· GWLR
· GWNBR
· Variance Inflation Factor
· Breusch-Pagan Test

Map Visualization
Map Visualization can help users to understand the circumstances of the observation area.

# RGUI using Shiny: FAST

# RGUI using C#: Wires

- For Spatial Data Analysis
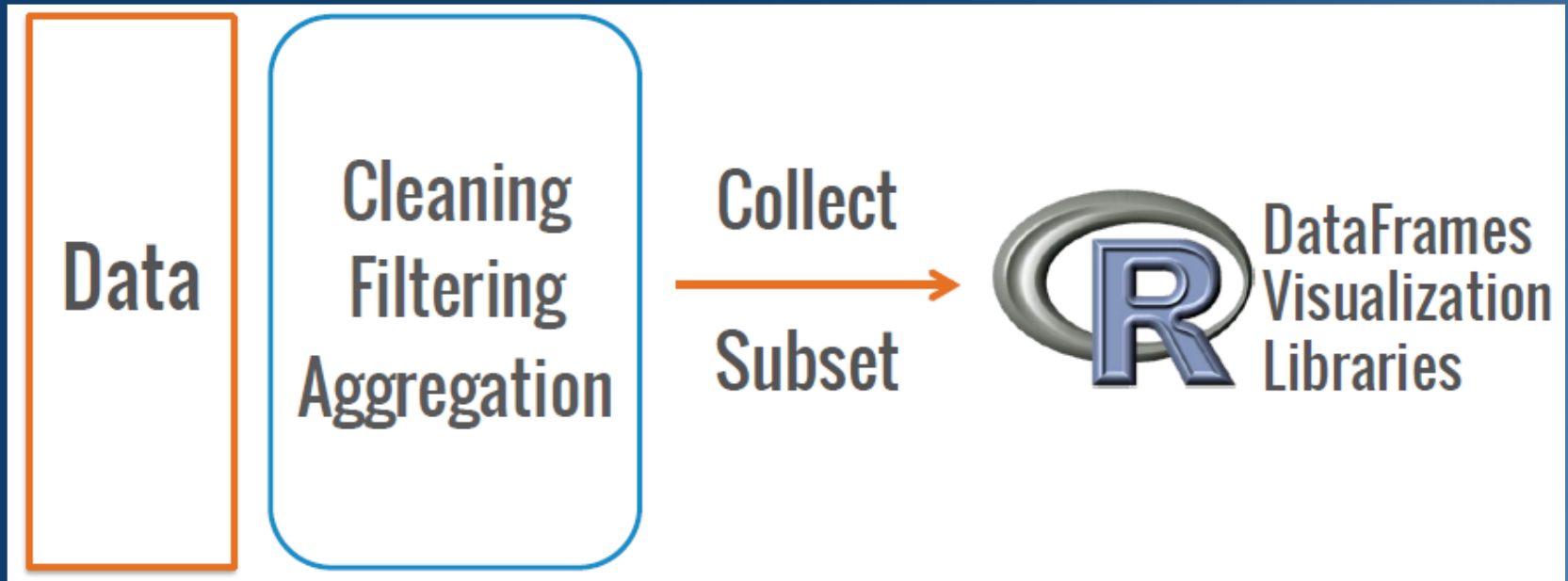
# RGUI using C#: Wires

# R Expert Modeler

# Dynamic Report Generation

Produce documents automatically: pdf, doc, html

Packages:

- Sweave
- knitr
- Markdown

# Big Data and R
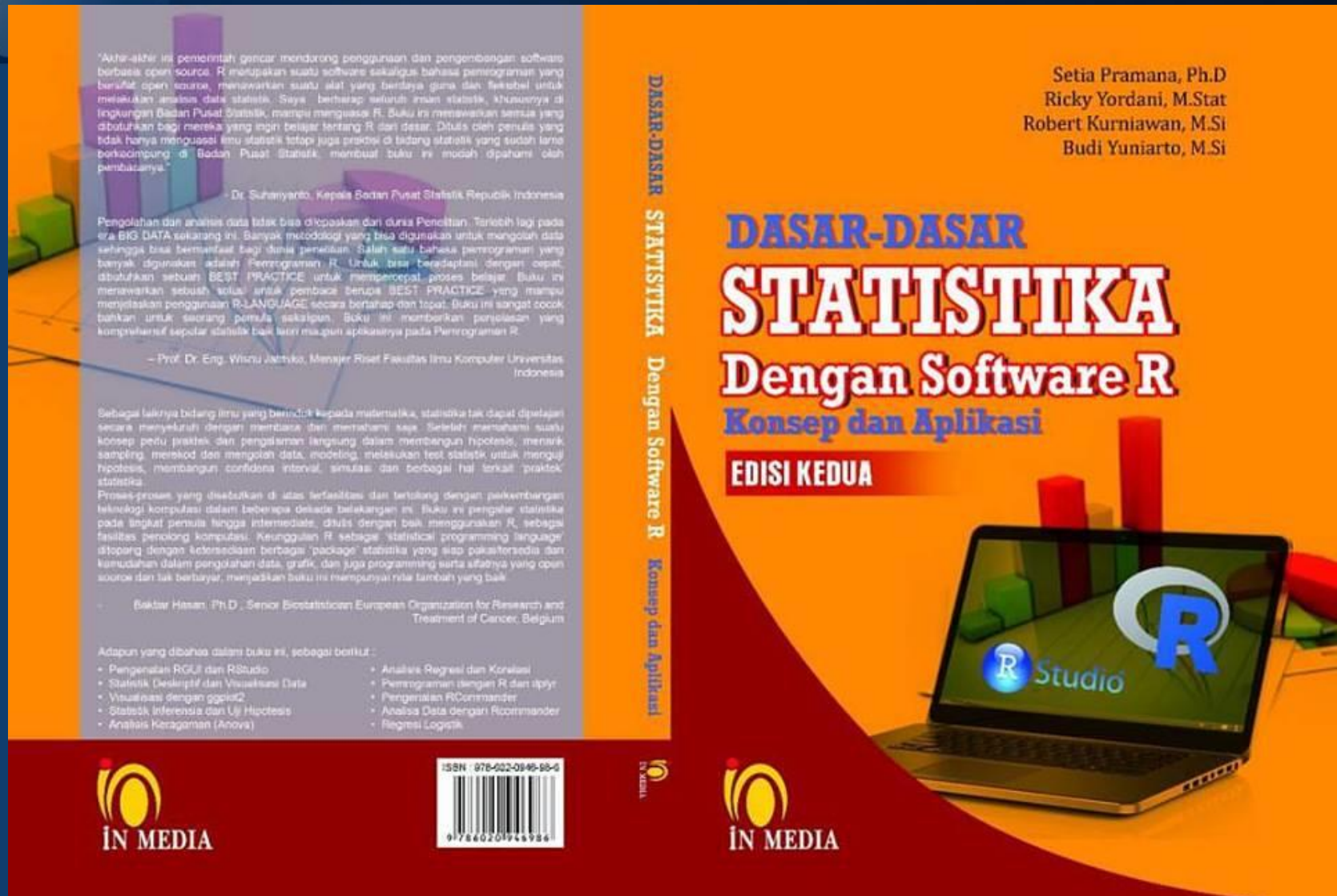
# Big Data

- **SparkR**



| Fast ! | | Statistics ! |
| Scalable | Spark + R | Packages |
| Interactive Shell | | Plots |

SparkR is a language binding that seamlessly integrates R with Spark, and enables native R programs to scale in a distributed setting

# A Start

# DATA MINING R dengan R

## Konsep Serta Implementasi

Kebutuhan akan eksplorasi dan analisis data semakin meningkat beberapa tahun terakhir. Metode eskplorasi dan analisis data juga mulai bergeser ke arah penggunaan data mining dan beberapa algoritma machine learning. Hal ini mendorong perubahan kurikulum dan materi yang harus disampaikan dan dikuasai mahasiswa khususnya mahasiswa jurusan statistik. Buku ini sangat saya rekomendasikan baik kepada mahasiswa maupun para pengajar karena buku ini tidak hanya memberikan teori namun juga mengajarkan bagaimana mengaplikasikan teori tersebut dalam contoh-contoh praktis. Buku ini juga memberikan keberagaman aplikasi dari data mining dengan tipe data yang berbeda-beda yang dapat diaplikasikan dengan software R.

**Dr. Erni Tri Astuti, M.Math - Direktur Politeknik Statistika STIS**

R merupakan salah satu alat pengolahan data yang sangat ampuh. Dengan bahasa yang lugas dan "to-the-point", penulis berhasil menyajikan data mining dengan pendekatan praktis menggunakan R. Buku ini merupakan batu pijakan yang sangat berguna buat para aspiring data scientist yang ingin menggeluti bidang data science

**Syafri Bahar S.Si., M.Sc., FRM - Vice President of Data Science GOJEK.**

Bahasan buku ini mencakup:

1. Pengantar Data Mining
2. Eksplorasi dan Visualisasi Data
3. Regresi Linear dan Logistik
4. Analisis Komponen Utama
5. Multivariate Anova
6. Supervised Learning (KNN, Decision Tree, Random Forest, dll)
7. Unsupervised Learning (Cluster Analysis)
8. Text Mining
9. Analisis Sentimen
10. Data Mining dalam Bioinformatika

**IN MEDIA**

UMUM
ISBN 978-602-6469-
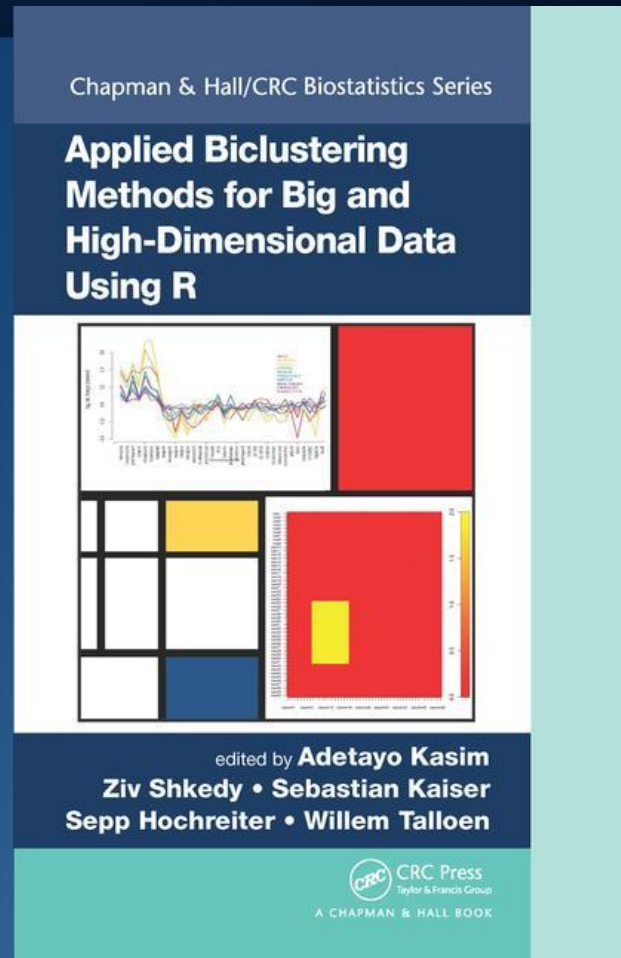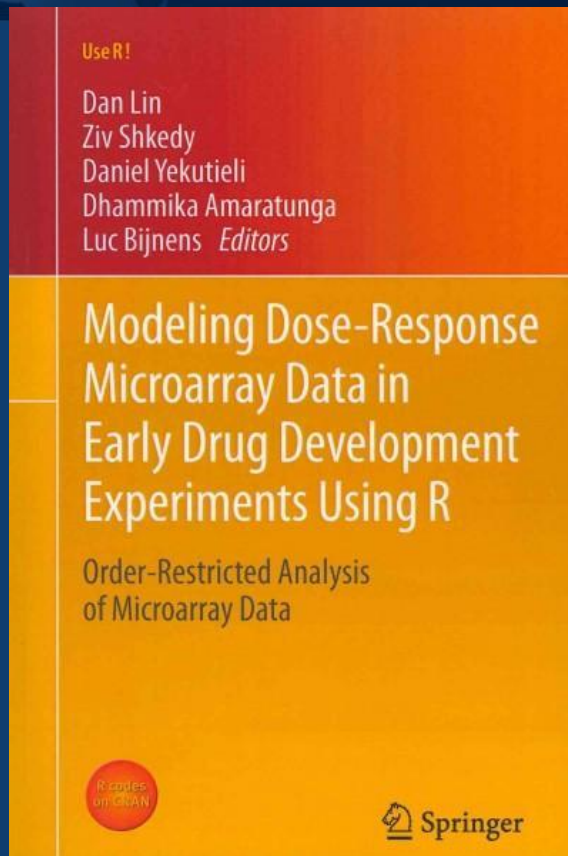
Harga P. Jawa Rp.

Setia Pramana, dkk

DATA MINING dengan R

Setia Pramana
Budi Yuniarto
Siti Mariyah
Ibnu Santoso
Rani Nooraeni

# DATA MINING R dengan R

## Konsep Serta Implementasi

**IN MEDIA**

8

# Book Chapters





Setia Pramana

# Conclusion

# Thank you for your attention!