
Version
0.3-0



*Programming with **B**ig **D**ata in **R***

Speaking Serial R with a Parallel Accent

Package Examples and Demonstrations

SPEAKING SERIAL R WITH A PARALLEL ACCENT (VER. 0.3-0)

pbdR PACKAGE EXAMPLES AND DEMONSTRATIONS

SEPTEMBER 25, 2015

DREW SCHMIDT
Business Analytics and Statistics
University of Tennessee

WEI-CHEN CHEN
pbdR Core Team

GEORGE OSTROUCHOV
Computer Science and Mathematics Division
Oak Ridge National Laboratory

PRAGNESHKUMAR PATEL
National Institute for Computational Sciences
University of Tennessee



VERSION 0.3-0

© 2012–2015 **p****b****d****R** Core Team. All rights reserved.

Permission is granted to make and distribute verbatim copies of this vignette and its source provided the copyright notice and this permission notice are preserved on all copies.

This manual may be incorrect or out-of-date. The authors assume no responsibility for errors or omissions, or for damages resulting from the use of the information contained herein.

Cover art from *Join, or Die* ([Franklin, 1754](#)). Illustrations were created using the **ggplot2** package ([Wickham, 2009](#)), native R functions, and Microsoft Powerpoint.

This publication was typeset using L^AT_EX.

Contents

| | |
|--|-----------|
| List of Figures | v |
| List of Tables | vii |
| Acknowledgements | viii |
| Note About the Cover | ix |
| Disclaimer | 1 |
| | |
| I Preliminaries | 2 |
| | |
| 1 Introduction | 3 |
| 1.1 What is pbdR? | 3 |
| 1.2 Why Parallelism? Why pbdR? | 5 |
| 1.3 Installation | 5 |
| 1.4 Structure of pbdDEMO | 6 |
| 1.4.1 List of Demos | 6 |
| 1.4.2 List of Benchmarks | 8 |
| 1.5 Exercises | 9 |
| | |
| 2 Background | 10 |
| 2.1 Parallelism | 10 |
| 2.2 SPMD Programming with R | 14 |
| 2.3 Notation | 15 |
| 2.4 Exercises | 15 |
| | |
| II Direct MPI Methods | 17 |
| | |
| 3 MPI for the R User | 18 |
| 3.1 MPI Basics | 18 |
| 3.2 pbdMPI vs Rmpi | 19 |
| 3.3 The GBD Data Structure | 21 |

| | | |
|------------|--|-----------|
| 3.4 | Common MPI Operations | 23 |
| 3.4.1 | Basic Communicator Wrangling | 24 |
| 3.4.2 | Reduce, Broadcast, and Gather | 25 |
| 3.4.3 | Printing and RNG Seeds | 26 |
| 3.4.4 | Apply, Lapply, and Sapply | 28 |
| 3.5 | Miscellaneous Basic MPI Tasks | 29 |
| 3.5.1 | Timing MPI Tasks | 29 |
| 3.5.2 | Distributed Logic | 29 |
| 3.6 | Exercises | 31 |
| 4 | Basic Statistics Examples | 34 |
| 4.1 | Monte Carlo Simulation | 34 |
| 4.2 | Sample Mean and Sample Variance | 37 |
| 4.3 | Binning | 38 |
| 4.4 | Quantile | 38 |
| 4.5 | Ordinary Least Squares | 40 |
| 4.6 | Exercises | 41 |
| III | Distributed Matrix Methods | 43 |
| 5 | DMAT | 44 |
| 5.1 | Block Data Distributions | 47 |
| 5.2 | Cyclic Data Distributions | 48 |
| 5.3 | Block-Cyclic Data Distributions | 49 |
| 5.4 | Summary | 52 |
| 5.5 | Exercises | 53 |
| 6 | Constructing Distributed Matrices | 54 |
| 6.1 | Fixed Global Dimension | 54 |
| 6.1.1 | Constructing Simple Distributed Matrices | 55 |
| 6.1.2 | Diagonal Distributed Matrices | 56 |
| 6.1.3 | Random Matrices | 57 |
| 6.2 | Fixed Local Dimension | 59 |
| 6.3 | Exercises | 59 |
| 7 | Basic Examples | 60 |
| 7.1 | Reductions and Transformations | 61 |
| 7.1.1 | Reductions | 61 |
| 7.1.2 | Transformations | 62 |
| 7.2 | Matrix Arithmetic | 62 |
| 7.3 | Matrix Factorizations | 63 |
| 7.4 | Exercises | 64 |
| 8 | Advanced Statistics Examples | 66 |
| 8.1 | Sample Mean and Variance Revisited | 66 |
| 8.2 | Verification of Distributed System Solving | 67 |
| 8.3 | Compression with Principal Components Analysis | 68 |

| | | |
|-----------|---|------------|
| 8.4 | Predictions with Linear Regression | 69 |
| 8.5 | Exercises | 70 |
| IV | Reading and Managing Data | 71 |
| 9 | Readers | 72 |
| 9.1 | CSV Files | 72 |
| 9.2 | Exercises | 73 |
| 10 | Parallel NetCDF4 Files | 74 |
| 10.1 | Introduction | 74 |
| 10.2 | Parallel Write and Read | 76 |
| 10.3 | Exercises | 78 |
| 11 | Redistribution Methods | 79 |
| 11.1 | Distributed Matrix Redistributions | 79 |
| 11.2 | Implicit Redistributions | 81 |
| 11.3 | Load Balance and Unload Balance | 82 |
| 11.4 | Convert Between GBD and DMAT | 84 |
| 11.5 | Exercises | 85 |
| V | Applications | 87 |
| 12 | Likelihood | 88 |
| 12.1 | Introduction | 88 |
| 12.2 | Normal Distribution | 89 |
| 12.3 | Likelihood Ratio Test | 90 |
| 12.4 | Multivariate Normal Distribution | 91 |
| 12.5 | Exercises | 92 |
| 13 | Model-Based Clustering | 93 |
| 13.1 | Introduction | 93 |
| 13.2 | Parallel Model-Based Clustering | 94 |
| 13.3 | An Example Using the <i>Iris</i> Dataset | 95 |
| 13.3.1 | <i>Iris</i> in Serial Code and Sample Outputs | 97 |
| 13.3.2 | <i>Iris</i> in GBD Code | 99 |
| 13.3.3 | <i>Iris</i> in <code>ddmatrix</code> Code | 100 |
| 13.4 | Exercises | 102 |
| 14 | Phylogenetic Clustering (Phyloclustering) | 103 |
| 14.1 | Introduction | 103 |
| 14.2 | The phyclust Package | 105 |
| 14.3 | Bootstrap Method | 106 |
| 14.4 | Task Pull Parallelism | 107 |
| 14.5 | An Example Using the <i>Pony 524</i> Dataset | 109 |
| 14.6 | Exercises | 110 |

| | |
|---|------------|
| 15 Bayesian MCMC | 111 |
| 15.1 Introduction | 111 |
| 15.2 Hastings-Metropolis Algorithm | 112 |
| 15.3 Galaxy Velocity | 114 |
| 15.4 Parallel Random Number Generator | 115 |
| 15.5 Exercises | 116 |
| 16 Pairwise Distance and Comparisons | 118 |
| 16.1 Introduction | 118 |
| 16.2 Distributed Distance and Comparisons | 119 |
| 16.3 Hierarchical Clustering | 120 |
| 16.4 Neighbor Joining | 121 |
| 16.5 Exercises | 122 |
| VI Appendix | 123 |
| A Numerical Linear Algebra and Linear Least Squares Problems | 124 |
| A.1 Forming the Normal Equations | 124 |
| A.2 Using the QR Factorization | 125 |
| A.3 Using the Singular Value Decomposition | 126 |
| B Linear Regression and Rank Degeneracy in R | 127 |
| VII Miscellany | 129 |
| References | 130 |
| Index | 135 |

List of Figures

| | | |
|------|--|-----|
| 1.1 | pbdr Packages | 4 |
| 1.2 | pbdr Package Use | 4 |
| 1.3 | pbdr Interface to Foreign Libraries | 5 |
| 2.1 | Task Parallelism Example | 11 |
| 2.2 | Task Parallelism Example | 13 |
| 4.1 | Approximating π | 36 |
| 5.1 | Matrix Distribution Schemes | 44 |
| 5.2 | Matrix Distribution Schemes Onto a 2-Dimensional Grid | 45 |
| 7.1 | Covariance Benchmark | 63 |
| 10.1 | Monthly averaged temperature | 76 |
| 11.1 | Matrix Redistribution Functions | 80 |
| 11.2 | Load Balancing/Unbalancing Data | 83 |
| 11.3 | Converting Between GBD and DMAT | 85 |
| 13.1 | Iris pair-wised scatter plot | 96 |
| 13.2 | Iris Clustering Plots — Serial | 99 |
| 13.3 | Iris Clustering Plots — GBD | 100 |
| 13.4 | Iris Clustering Plots — GBD | 101 |
| 14.1 | Retrovirus phylogeny originated from Weiss (2006). | 105 |
| 14.2 | 146 EIAV sequences of <i>Pony 524</i> in three clusters. | 106 |
| 15.1 | Histograms of velocities of 82 galaxies | 114 |
| 15.2 | MCMC results of velocities of 82 galaxies | 115 |
| 16.1 | Hierarchical clustering result of <i>iris</i> dataset. | 120 |

| | |
|--|-----|
| 16.2 Neighbor-joining tree of <i>Pony 524</i> dataset colored by three clusters. | 121 |
|--|-----|

List of Tables

| | | |
|------|---|----|
| 3.1 | Benchmark Comparing Rmpi and pbdMPI | 21 |
| 5.1 | Processor Grid Shapes with 6 Processors | 46 |
| 10.1 | Functions for accessing NetCDF4 files | 75 |
| 11.1 | Implicit Data Redistributions | 81 |
| 13.1 | Parallel Mode-Based Clustering Algorithms in pmclust | 95 |

Acknowledgements

Schmidt, Ostrouchov, and Patel were supported in part by the project “NICS Remote Data Analysis and Visualization Center” funded by the Office of Cyberinfrastructure of the U.S. National Science Foundation under Award No. ARRA-NSF-OCI-0906324 for NICS-RDAV center.

Chen was supported in part by the project “Bayesian Assessment of Safety Profiles for Pregnant Women From Animal Study to Human Clinical Trial” funded by U.S. Food and Drug Administration, Office of Women’s Health. The project was supported in part by an appointment to the Research Participation Program at the Center For Biologics Evaluation and Research administered by the Oak Ridge Institute for Science and Education through an interagency agreement between the U.S. Department of Energy and the U.S. Food and Drug Administration

Chen was supported in part by the Department of Ecology and Evolutionary Biology at the University of Tennessee, Knoxville, and a grant from the National Science Foundation (MCB-1120370.)

Chen and Ostrouchov were supported in part by the project “Visual Data Exploration and Analysis of Ultra-large Climate Data” funded by U.S. DOE Office of Science under Contract No. DE-AC05-00OR22725.

This work used resources of National Institute for Computational Sciences at the University of Tennessee, Knoxville, which is supported by the Office of Cyberinfrastructure of the U.S. National Science Foundation under Award No. ARRA-NSF-OCI-0906324 for NICS-RDAV center. This work also used resources of the Oak Ridge Leadership Computing Facility at the Oak Ridge National Laboratory, which is supported by the Office of Science of the U.S. Department of Energy under Contract No. DE-AC05-00OR22725. This work used resources of the Newton HPC Program at the University of Tennessee, Knoxville.

We also thank Brian D. Ripley, Kurt Hornik, Uwe Ligges, and Simon Urbanek from the R Core Team for discussing package release issues and helping us solve portability problems on different platforms.

Note About the Cover

The picture on the cover, *Join, or Die* (Franklin, 1754), is from a well-known political cartoon to those familiar with United States history. It has since become known as a symbol of unity to Americans. For more historical information and context, the reader is encouraged to see the well-documented Wikipedia entry https://en.wikipedia.org/wiki/Join%2C_or_Die#Legacy_of_the_cartoon.

The inclusion of this symbol here is not intended to be political. In the context of distributed computing, the multiple nodes must "join or die". In this document, we will learn how to join together multiple nodes to solve problems in statistics using MPI, accessible from R.

Disclaimer

Warning: The findings and conclusions in this article have not been formally disseminated by the U.S. Department of Health & Human Services nor by the U.S. Department of Energy, and should not be construed to represent any determination or policy of University, Agency, Administration and National Laboratory.

This document is written to explain the main functions of **pbdDEMO** (Schmidt et al., 2013), version 0.3-0. Every effort will be made to ensure future versions are consistent with these instructions, but features in later versions may not be explained in this document.

Information about the functionality of this package, and any changes in future versions can be found on website: “Programming with Big Data in R” at <http://r-pbd.org/>.

Part I

Preliminaries

There are things which seem incredible to most men who have not studied Mathematics.

—Archimedes of Syracuse

1.1 What is pbdR?

The “Programming with Big Data in R” project (Ostrouchov et al., 2012) (**pbdR** for short) is a project that aims to elevate the statistical programming language R (R Core Team, 2012a) to leadership-class computing platforms. The main goal is empower data scientists by bringing flexibility and a big analytics toolbox to big data challenges, with an emphasis on productivity, portability, and performance. We achieve this in part by mapping high-level programming syntax to portable, high-performance, scalable, parallel libraries. In short, we make R scalable.

Figure 1.1 shows the current list of **pbdR** packages released to the CRAN (<http://cran.r-project.org>), and how they depend on each other. More explicitly, the current **pbdR** packages (Chen et al., 2012a,c; Schmidt et al., 2012a,c; Patel et al., 2013a; Schmidt et al., 2013) are:

- **pbdMPI** — an efficient interface to MPI (Gropp et al., 1994) with a focus on Single Program/Multiple Data (SPMD) parallel programming style.
- **pbdSLAP** — bundles scalable dense linear algebra libraries in double precision for R, based on ScaLAPACK version 2.0.2 (Blackford et al., 1997).
- **pbdNCDF4** — interface to Parallel Unidata NetCDF4 format data files (NetCDF Group, 2008).
- **SEXPtools** — SEXP tools (?).
- **pbdBASE** — low-level ScaLAPACK codes and wrappers.
- **pbdDMAT** — distributed matrix classes and computational methods, with a focus on linear algebra and statistics.

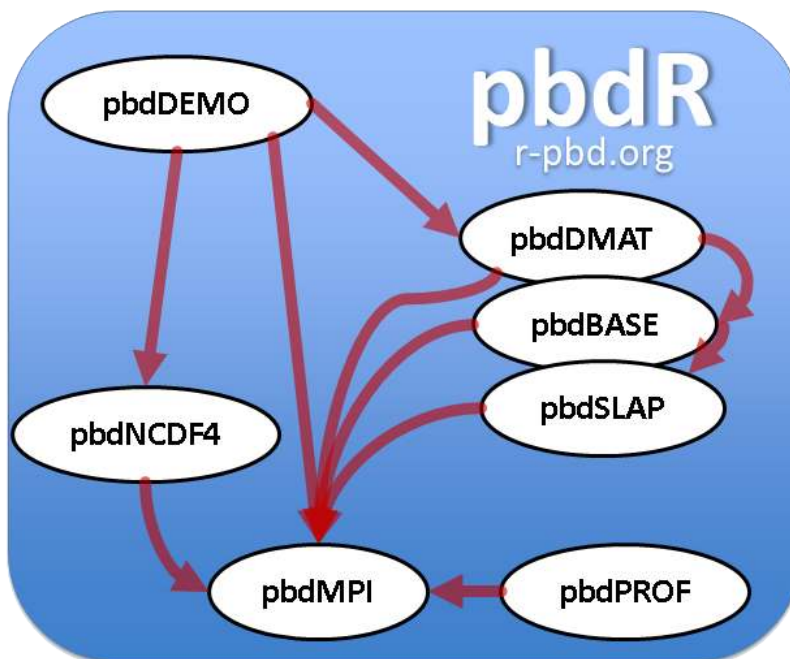


Figure 1.1: pbdR Packages

- **pbdDEMO** — set of package demonstrations and examples, and this unifying vignette.

To try to make this landscape a bit more clear, one could divide **pbdR** packages into those meant for users, developers, or something in-between. Figure 1.2 shows a gradient scale representation,

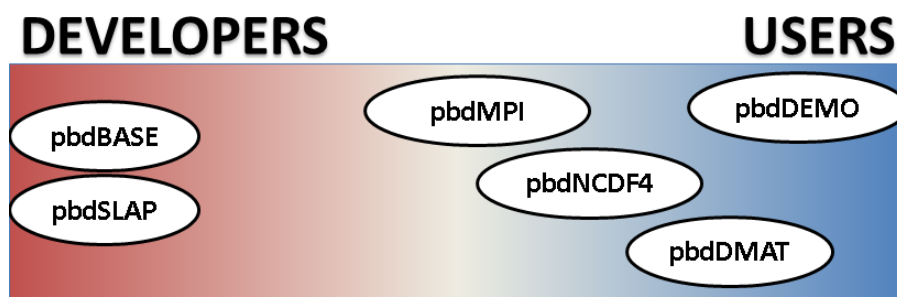


Figure 1.2: pbdR Package Use

where more red means the package is more for developers, while more blue means the package is more for users. For example, **pbdDEMO** is squarely meant for users of **pbdR** packages, while **pbdBASE** and **pbdSLAP** are really not meant for general use. The other packages fall somewhere in-between, having plenty of utility for both camps.

Finally, Figure 1.3 shows **pbdR** relationship to high-performance libraries.

In this vignette, we offer many examples using the above **pbdR** packages. Many of the examples are high-level applications and may be commonly found in basic Statistics. The purpose is to show how to reuse the preexisting functions and utilities of **pbdR** to create minor extensions which can quickly solve problems in an efficient way. The reader is encouraged to reuse and

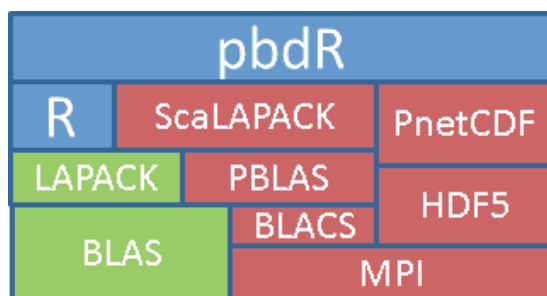


Figure 1.3: pbdR Interface to Foreign Libraries

re-purpose these functions.

The **pbdDEMO** package consists of two main parts. The first is a collection of roughly 20+ package demos. These offer example uses of the various **pbdR** packages. The second is this vignette, which attempts to offer detailed explanations for the demos, as well as sometimes providing some mathematical or statistical insight. A list of all of the package demos can be found in Section 1.4.1.

1.2 Why Parallelism? Why pbdR?

It is common, in a document such as this, to justify the need for parallelism. Generally this process goes:

Blah blah blah Moore's Law, blah blah Big Data, blah blah blah Concurrency.

How about this? Parallelism is cool. Any boring nerd can use one computer, but using 10,000 at once is another story. We don't call them *supercomputers* for nothing.

But unfortunately, lots of people who would otherwise be thrilled to do all kinds of cool stuff with massive behemoths of computation — computers with names like **KRAKEN**¹ and **TITAN**² — are burdened by an unfortunate reality: it's really, really hard. Enter **pbdR**. Through our project, we put a shiny new set of clothes on high-powered compiled code, making massive-scale computation accessible to a wider audience of data scientists than ever before. Analytics in supercomputing shouldn't just be for the elites.

1.3 Installation

One can download **pbdDEMO** from CRAN at <http://cran.r-project.org>, and the installation can be done with the following commands

```
tar zxvf pbdDEMO_0.2-0.tar.gz
R CMD INSTALL pbdDEMO
```

¹ <http://www.nics.tennessee.edu/computing-resources/kraken>

² <http://www.olcf.ornl.gov/titan/>

Since **pbdEMO** depends on other **pbdR** packages, please read the corresponding vignettes if installation of any of them is unsuccessful.

1.4 Structure of pbdDEMO

The **pbdDEMO** package consists of several key components:

1. This vignette
2. A set of demos in the `demo/` tree
3. A set of benchmark codes in the `Benchmarks/` tree

The following subsections elaborate on the contents of the latter two.

1.4.1 List of Demos

A full list of demos contained in the **pbdDEMO** package is provided below. We may or may not describe all of the demos in this vignette.

List of Demos

```
### (Use Rscript.exe for windows systems)

# ----- #
# II Direct MPI Methods #
# ----- #

### Chapter 4
# Monte carlo simulation
mpiexec -np 4 Rscript -e "demo(monte_carlo, package='pbdDMAT', ask=F,
  echo=F)"
# Sample mean and variance
mpiexec -np 4 Rscript -e "demo(sample_stat, package='pbdDMAT', ask=F,
  echo=F)"
# Binning
mpiexec -np 4 Rscript -e "demo(binning, package='pbdDMAT', ask=F,
  echo=F)"
# Quantile
mpiexec -np 4 Rscript -e "demo(quantile, package='pbdDMAT', ask=F,
  echo=F)"
# OLS
mpiexec -np 4 Rscript -e "demo(ols, package='pbdDMAT', ask=F, echo=F)"
# Distributed Logic
mpiexec -np 4 Rscript -e "demo(comparators, package='pbdDMAT', ask=F,
  echo=F)"

# ----- #
# III Distributed Matrix Methods #
# ----- #
```

```

### Chapter 6
# Random matrix generation
mpiexec -np 4 Rscript -e "demo(randmat_global, package='pbdDMAT',
    ask=F, echo=F)"
mpiexec -np 4 Rscript -e "demo(randmat_local, package='pbdDMAT', ask=F,
    echo=F)"

### Chapter 8
# Sample statistics revisited
mpiexec -np 4 Rscript -e "demo(sample_stat_dmat, package='pbdDMAT',
    ask=F, echo=F)"
# Verify solving  $Ax=b$  at scale
mpiexec -np 4 Rscript -e "demo(verify, package='pbdDMAT', ask=F,
    echo=F)"
# PCA compression
mpiexec -np 4 Rscript -e "demo(pca, package='pbdDMAT', ask=F, echo=F)"
# OLS and predictions
mpiexec -np 4 Rscript -e "demo(ols_dmat, package='pbdDMAT', ask=F,
    echo=F)"

# ----- #
# IV Reading and Managing Data #
# ----- #

### Chapter 9
mpiexec -np 4 Rscript -e "demo(read_csv, package='pbdDMAT', ask=F,
    echo=F)"

### Chapter 10
# Reading and writing parallel NetCDF4
Rscript -e "demo(trefht, package="pbdDEMO", ask = F, echo = F)"
mpiexec -np 4 Rscript -e "demo(nc4_serial, package='pbdDEMO', ask=F,
    echo=F)"
mpiexec -np 4 Rscript -e "demo(nc4_parallel, package='pbdDEMO', ask=F,
    echo=F)"
mpiexec -np 4 Rscript -e "demo(nc4_dmat, package='pbdDEMO', ask=F,
    echo=F)"
mpiexec -np 4 Rscript -e "demo(nc4_gbdc, package='pbdDEMO', ask=F,
    echo=F)"
mpiexec -np 4 Rscript -e "demo(nc4_par_write_1d, package='pbdDEMO',
    ask=F, echo=F)"
mpiexec -np 4 Rscript -e "demo(nc4_par_write_2d, package='pbdDEMO',
    ask=F, echo=F)"
mpiexec -np 4 Rscript -e "demo(nc4_par_read_1d, package='pbdDEMO',
    ask=F, echo=F)"
mpiexec -np 4 Rscript -e "demo(nc4_par_read_2d, package='pbdDEMO',
    ask=F, echo=F)"

### Chapter 11
# Load/unload balance
mpiexec -np 4 Rscript -e "demo(balance, package='pbdDMAT', ask=F,
    echo=F)"

```

```

# GBD to DMAT
mpiexec -np 4 Rscript -e "demo(gbd_dmat, package='pbdDMAT', ask=F,
    echo=F)"
# Distributed matrix redistributions
mpiexec -np 4 Rscript -e "demo(reblock, package='pbdDMAT', ask=F,
    echo=F)"

# ----- #
# V Applications #
# ----- #

### Chapter 13
# Parallel Model-Based Clustering
Rscript -e "demo(iris_overlap, package='pbdDEMO', ask=F, echo=F)"
Rscript -e "demo(iris_serial, package='pbdDEMO', ask=F, echo=F)"
Rscript -e "demo(iris_gbd, package='pbdDEMO', ask=F, echo=F)"
Rscript -e "demo(iris_dmat, package='pbdDEMO', ask=F, echo=F)"

### Chapter 14
mpiexec -np 4 Rscript -e "demo(task_pull, package='pbdMPI', ask=F,
    echo=F)"
mpiexec -np 4 Rscript -e "demo(phyclust_bootstrap, package='pbdDEMO',
    ask=F, echo=F)"

### Chapter 15
mpiexec -np 4 Rscript -e "demo(mcmc_galaxy, package='pbdDEMO', ask=F,
    echo=F)"

### Chapter 16
mpiexec -np 4 Rscript -e "demo(dist_iris, package='pbdDEMO', ask=F,
    echo=F)"
mpiexec -np 4 Rscript -e "demo(dist_pony, package='pbdDEMO', ask=F,
    echo=F)"

```

1.4.2 List of Benchmarks

At the time of writing, there are benchmarks for computing covariance, linear models, and principal components. The benchmarks come in two variants. The first is an ordinary set of benchmark codes, which generate data of specified dimension(s) and time the indicated computation. This operation is replicated for a user-specified number of times (default 10), and then the timing results are printed to the terminal.

From the **Benchmarks/** subtree, the user can run the first set of benchmarks with, for example, 4 processors by issuing any of the commands:

```

### (Use Rscript.exe for windows systems)
mpiexec -np 4 Rscript cov.r
mpiexec -np 4 Rscript lmfit.r
mpiexec -np 4 Rscript pca.r

```

The second set of benchmarks are those that try to find the “balancing” point where, for the indicted computation with user specified number of cores, the computation is performed faster using **pbdR** than using serial R. In general, throwing a bunch of cores at a problem may not be the best course of action, because parallel algorithms (almost always) have inherent overhead over their serial counterparts that can make their use ill-advised for small problems. But for sufficiently big (which is usually not very big at all) problems, that overhead should quickly be dwarfed by the increased scalability.

From the **Benchmarks/** subtree, the user can run the second set of benchmarks with, for example, 4 processors by issuing any of the commands:

```
### (Use Rscript.exe for windows systems)
mpiexec -np 4 Rscript balance_cov.r
mpiexec -np 4 Rscript balance_lmfit.r
mpiexec -np 4 Rscript balance_pca.r
```

Now we must note that there are other costs than just statistical computation. There is of course the cost of disk IO (when dealing with real data). However, a parallel file system should help with this, and for large datasets should actually be faster anyway. The main cost not measured here is the cost of starting all of the R processes and loading packages. Assuming R is not compiled statically (and it almost certainly is not), then this cost is non-trivial and somewhat unique to very large scale computing. For instance, it took us well over an hour to start 12,000 R sessions and load the required packages on the supercomputer KRAKEN³. This problem is not unique to R, however. It affects any project that has a great deal of dynamic library loading to do. This includes Python, although their community has made some impressive strides in dealing with this problem.

1.5 Exercises

- 1-1 Read the MPI wikipedia page https://en.wikipedia.org/wiki/Message_Passing_Interface including it's history, overview, functionality, and concepts sections.
- 1-2 Read the **pbdMPI** vignette and install either OpenMPI (<http://www.open-mpi.org/>) or MPICH2 (<http://www.mcs.anl.gov/research/projects/mpich2/>), and test if the installation is correct (see <http://www.r-pbd.org/install.html> for more details).
- 1-3 After completing Exercise 1-2, install all **pbdR** packages and run each package's demo codes.

³See [https://en.wikipedia.org/wiki/Kraken_\(supercomputer\)](https://en.wikipedia.org/wiki/Kraken_(supercomputer))

*We stand at the threshold of a many core world.
The hardware community is ready to cross this
threshold. The parallel software community is
not.*

—Tim Mattson

2.1 Parallelism

What is parallelism? At its core (pun intended), parallelism is all about trying to throw more resources at a problem, usually to get the problem to complete faster than it would with the more minimal resources. Sometimes we wish to utilize more resources as a means of being able to make a computationally (literally or practically) intractable problem into one which will complete in a reasonable amount of time. Somewhat more precisely, parallelism is the leveraging of parallel processing. It is a general programming model whereby we execute different computations simultaneously. This stands in contrast to *serial* programming, where you have a stream of commands, executed one at a time.

Serial programming has been the dominant model from the invention of the computer to present, although this is quickly changing. The reasons why this is changing are numerous and boring; the fact is, if it is true now that a researcher must know some level of programming to do his/her job, then it is certainly true that in the near future that he/she will have to be able to do some parallel programming. Anyone who would deny this is, frankly, more likely trying to vocally assuage personal fears more so than accurately forecasting based on empirical evidence. For many, parallel programming isn't *coming*; it's *here*.

As a general rule, parallelism should only come after you have exhausted serial optimization. Even the most complicated parallel programs are made up of serial pieces, so inefficient serial codes produce inefficient parallel codes. Also, generally speaking, one can often eke out much better performance by implementing a very efficient serial algorithm rather than using a handful of cores (like on a modern multicore laptop) using an inefficient parallel algorithm. However, once that serial-optimization well runs dry, if you want to continue seeing performance gains,

then you must implement your code in parallel.

Next, we will discuss some of the major parallel programming models. This discussion will be fairly abstract and superficial; however, the overwhelming bulk of this text is comprised of examples which will appeal to data scientists, so for more substantive examples, especially for those more familiar with parallel programming, you may wish to jump to Section 4.

Data Parallelism

There are many ways to write parallel programs. Often these will depend on the physical hardware you have available to you (multicore laptop, several GPU's, a distributed supercomputer, ...). The `pbdr` project is principally concerned with *data parallelism*. We will expand on the specifics in Section 2.3 and provide numerous examples throughout this guide. However, in general, data parallelism is a parallel programming model whereby the programmer splits up a data set and applies operations on the sub-pieces to solve one larger problem.

Figure 2.1 offers a visualization of a very simple data parallelism problem. Say we have an array

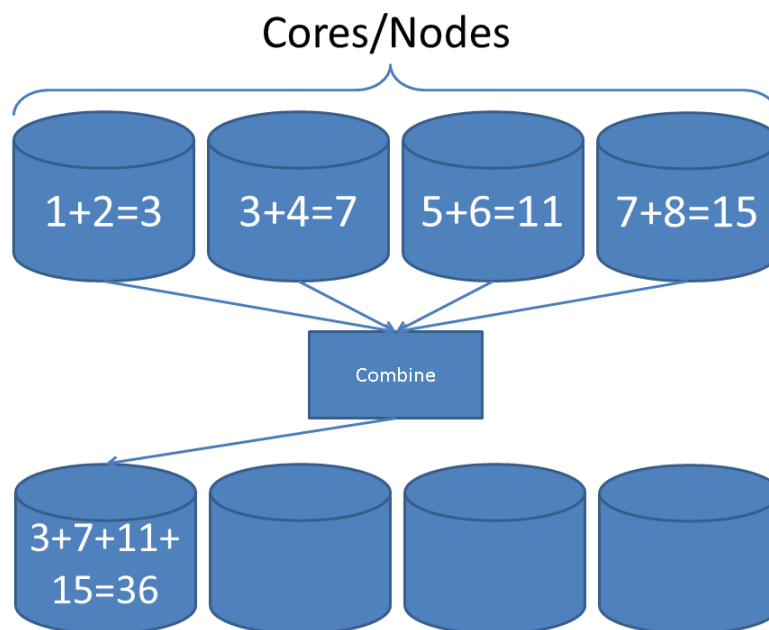


Figure 2.1: Task Parallelism Example

consisting of the values 1 through 8, and we have 4 cores (processing units) at our disposal, and we want to add up all of the elements of this array. We might distribute the data as in the diagram (the first two elements of the array on the first core, the next two elements on the second core, and so on). We then perform a local summation operation; this local operation is serial, but because we have divided up the overall task of summation across the multiple processors, for a very large array we would expect to see performance gains.

A very loose pseudo code for this procedure might look like:

Pseudocode

```
1: mydata = map(data)
2: total_local = sum(mydata)
3: total = reduce(total_local)
4: if this_processor == processor_1 then
5:   print(total)
6: end if
```

Then each of the four cores could execute this code simultaneously, with some cooperation between the processors for step 1 (in determining who owns what) and for the reduction in step 3. This is an example of using a higher-level parallel programming paradigm called “Single Program/Multiple Data” or SPMD. . We will elucidate more as to exactly what this means in the sections to follow.

Task Parallelism

Data parallelism is one parallel programming model. By contrast, another important parallel programming model is *task parallelism*, which much of the R community is already fairly adept at, especially when using the manager/workers paradigm (more on this later). Common packages for this kind of work include **snow** (Tierney et al., 2012), **parallel** (R Core Team, 2012b), and **Rmpi** (Yu, 2002)¹.

Task parallelism involves, as the name implies, distributing different execution tasks across processors. Task parallelism is often *embarrassingly parallel* — meaning the parallelism is so easy to exploit that it is embarrassing. This kind of situation occurs when you have complete independence, or a *loosely coupled* problem (as opposed to something *tightly coupled*, like computing the singular value decomposition (SVD) of a distributed data matrix, for example).

As a simple example of task parallelism, say you have one dataset and four processing cores, and you want to fit all four different linear regression models for that dataset, and then choose the model with lowest AIC (Akaike, 1974) (we are not arguing that this is necessarily a good idea; this is just an example). Fitting one model does not have any dependence on fitting another, so you might want to just do the obvious thing and have each core fit a separate model, compute the AIC value locally, then compare all computed AIC values, lowest is the winner. Figure 2.2 offers a simple visualization of this procedure.

A very loose pseudo code for this problem might look like:

Pseudocode

¹For more examples, see “CRAN Task View: High-Performance and Parallel Computing with R” at <http://cran.r-project.org/web/views/HighPerformanceComputing.html>.

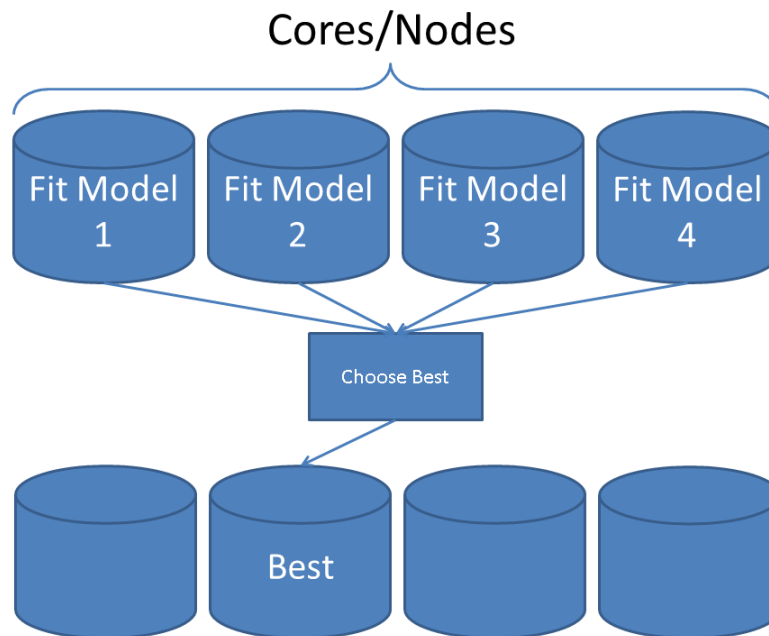


Figure 2.2: Task Parallelism Example

```

1: load_data()
2: if this_processor == processor_1 then
3:     distribute_tasks()
4: else
5:     receive_tasks()
6: end if
7: model_aic = aic(fit(mymodel))
8: best_aic = min(allgather(model_aic))
9: if model_aic == best_aic then
10:     print(mymodel)
11: end if

```

Then each of the four cores could execute this code simultaneously, with some cooperation between the processors for the distribution of tasks (which model to fit) in step 2 and for the gather operation in step 8.

The line between data parallelism and task parallelism sometimes blurs, especially in simple examples such as those presented here; given our loose definitions of terms, is our first example really data parallelism? Or is it task parallelism? It is best not to spend too much time worrying about what to call it and instead focus on how to do it. These are simple examples and should not be taken too far out of context. All that said, the proverbial rabbit hole of parallel programming goes quite deep, and often it is not a matter of one programming model or another, but leveraging several at once to solve a complicated problem.

2.2 SPMD Programming with R

Throughout this document, we will be using the ‘Single Program/Multiple Data’, or SPMD, paradigm for distributed computing. Writing programs in the SPMD style is a very natural way of doing computations in parallel, but can be somewhat difficult to properly describe. As the name implies, only one program is written, but the different processors involved in the computation all execute the code independently on different portions of the data. The process is arguably the most natural extension of running serial codes in batch. This model lends itself especially well to data parallelism problems.

Unfortunately, executing jobs in batch is a somewhat unknown way of doing business to the typical R user. While details and examples about this process will be provided in the chapters to follow, the reader is encouraged to read the **pbdMPI** package’s vignette (Chen et al., 2012b) first. Ideally, readers should run the demos of the **pbdMPI** package, going through the code step by step.

This paradigm is just one model of parallel programming, and in reality, is a sort of “meta model”, encompassing many other parallel programming models. The R community is already familiar with the manager/workers² programming model. This programming model is particularly well-suited to task parallelism, where generally one processor will distribute the task load to the other processors.

The two are not mutually exclusive, however. It is easy enough to engage in task parallelism from an SPMD-written program. To do this, you essentially create a “task-parallelism block”, where one processor

Pseudocode

```
1: if this_processor == manager then  
2:   distribute_tasks()  
3: else  
4:   receive_tasks()  
5: end if
```

See Exercise 2-2 for more details.

One other model of note is the MapReduce model. A well-known implementation of this is Apache’s Hadoop, which is really more of a poor man’s distributed file system with MapReduce bolted on top. The R community has a strange affection for MapReduce, even among people who have never used it. MapReduce, and for instance Hadoop, most certainly has its place, but one should be aware that MapReduce is not very well-suited for tightly coupled problems; this difficulty goes beyond the fact that tightly coupled problems are harder to parallelize than their embarrassingly parallel counterparts, and is, in part, inherent to MapReduce itself. For the remainder of this document, we will not discuss MapReduce further.

²Sometimes referred to as “master/slaves” or “master/workers”

2.3 Notation

Note that we tend to use suffix `.gbd` for an object when we wish to indicate that the object is “general block distributed.” This is purely for pedagogical convenience, and has no semantic meaning. Since the code is written in SPMD style, you can think of such objects as referring to either a large, global object, or to a processor’s local piece of the whole (depending on context). This is less confusing than it might first sound.

We will not use this suffix to denote a global object common to all processors. As a simple example, you could imagine having a large matrix with (global) dimensions $m \times n$ with each processor owning different collections of rows of the matrix. All processors might need to know the values for m and n ; however, m and n do not depend on the local process, and so these do not receive the `.gbd` suffix. In many cases, it may be a good idea to invent an S4 class object and a corresponding set of methods. Doing so can greatly improve the usability and readability of your code, but is never strictly necessary. However, these constructions are the foundation of the **pbdBASE** (Schmidt et al., 2012a) and **pbdDMAT** (Schmidt et al., 2012c) packages.

On that note, depending on your requirements in distributed computing with R, it may be beneficial to you to use higher **pbdR** toolchain. If you need to perform dense matrix operations, or statistical analysis which depend heavily on matrix algebra (linear modeling, principal components analysis, ...), then the **pbdBASE** and **pbdDMAT** packages are a natural choice. The major hurdle to using these tools is getting the data into the appropriate **ddmatrix** format, although we provide many tools to ease the pains of doing so. Learning how to use these packages can greatly improve code performance, and take your statistical modeling in R to previously unimaginable scales.

Again for the sake of understanding, we will at times append the suffix `.dmat` to objects of class **ddmatrix** to differentiate them from the more general `.gbd` object. As with `.gbd`, this carries no semantic meaning, and is merely used to improve the readability of example code (especially when managing both “`.gbd`” and **ddmatrix** objects).

2.4 Exercises

- 2-1 Read the SPMD wiki page at <http://en.wikipedia.org/wiki/SPMD> and it’s related information.
- 2-2 **pbdMPI** provides a function `get.jid()` to divide N jobs into n processors nearly equally which is best for homogeneous computing environment to do task parallelism. The FAQs section of **pbdMPI**’s vignette has an example, try it as next.

R Code

```
1 library(pbdMPI, quiet=TRUE)
2 init()
3
4 id <- get.jid(N)
5
```

```
6  ### Using a loop
7  for(i in id){
8      # put independent task i script here
9  }
10 finalize()
```

See Section [14.4](#) for more efficient task parallelism.

- 2-3 Multi-threading and forking are also popular methods of parallelism for shared memory systems, such as in a personal laptop. The function `mclapply()`³ in **parallel** originated from the **multicore** ([Urbanek, 2011](#)) package, and is for simple parallelism on shared memory machines by using the `fork` mechanism. Compare this with OpenMP ([OpenMP ARB, 1997](#)).

³This method is not available on Windows, because Windows has no system-level `fork` command.

Part II

Direct MPI Methods

Everybody who learns concurrency thinks they understand it, ends up finding mysterious races they thought weren't possible, and discovers that they didn't actually understand it yet after all.

—Herb Sutter

Cicero once said that “If you have a garden and a library, you have everything you need.” So in that spirit, for the next two chapters we will use the MPI library to get our hands dirty and root around in the dirt of low-level MPI programming.

3.1 MPI Basics

In a sense, Cicero (in the above tortured metaphor) was quite right. MPI is all that we *need* in the same way that I might only *need* bread and cheese, but really what I *want* is a pizza. MPI is somewhat low-level and can be quite fiddly, but mastering it adds a very powerful tool to the repertoire of the parallel R programmer, and is essential for anyone who wants to do large scale development of parallel codes.

“MPI” stands for “Message Passing Interface”. How it really works goes *well* beyond the scope of this document. But at a basic level, the idea is that the user is running a code on different compute nodes that (usually) can not directly modify objects in each others’ memory. In order to have all of the nodes working together on a common problem, data and computation directives are passed around over the network (often over a specialized link called infiniband).

At its core, MPI is a standard interface for managing communications (data and instructions) between different nodes or computers. There are several major implementations of this standard, and the one you should use may depend on the machine you are using. But this is a compiling issue, so user programs are unaffected beyond this minor hurdle. Some of the most well-known implementations are OpenMPI, MPICH2, and Cray MPT.

At the core of using MPI is the *communicator*. At a technical level, a communicator is a pretty complicated data structure, but these deep technical details are not necessary for proceeding. We will instead think of it somewhat like the post office. When we wish to send a letter (communication) to someone else (another processor), we merely drop the letter off at a post office mailbox (communicator) and trust that the post office (MPI) will handle things accordingly (sort of).

The general process for directly — or indirectly — utilizing MPI in SPMD programs goes something like this:

1. Initialize communicator(s).
2. Have each process read in its portion of the data.
3. Perform computations.
4. Communicate results.
5. Shut down the communicator(s).

Some of the above steps may be swept away under a layer of abstraction for the user, but the need may arise where directly interfacing with MPI is not only beneficial, but necessary.

More details and numerous examples using MPI with R are available in the sections to follow, as well as in the **pbdMPI** vignette.

3.2 pbdMPI vs Rmpi

There is another package on the CRAN which the R programmer may use to interface with MPI, namely **Rmpi** (Yu, 2002). There are several issues one must consider when choosing which package to use if one were to only use one of them.

1. (+) **pbdMPI** is easier to install than **Rmpi**
2. (+) **pbdMPI** is easier to use than **Rmpi**
3. (+) **pbdMPI** can often outperform **Rmpi**
4. (+) **pbdMPI** integrates with the rest of **pbdR**
5. (–) **Rmpi** can be used with **foreach** (Analytics, 2012) via **doMPI** (Weston, 2010)
6. (–) **Rmpi** can be used in the manager/worker paradigm

We do not believe that the above can be reduced to a zero-sum game with unambiguous winner and loser. Ultimately the needs of the user (or developer) are paramount. We believe that **pbdR** makes a very good case for itself, especially in the SPMD style, but it can not satisfy everyone. However, for the remainder of this section, we will present the case for several of the, as yet, unsubstantiated pluses above.

In the case of ease of use, **Rmpi** uses bindings very close to the level as they are used in C’s MPI API. Specifically, whenever performing, for example, a reduction operation such as “allreduce”, you must specify the type of your data. For example, using **Rmpi**’s API

```
1 mpi.allreduce(x, type = 1)
```

would perform the sum allreduce if the object **x** consists of integer data, while

```
1 mpi.allreduce(x, type = 2)
```

would be used if **x** consists of doubles. However, with **pbdMPI**

```
1 allreduce(x)
```

is used for both by making use of R’s S4 system of object oriented programming. This is not mere code golfing¹ that we are engaging in. The concept of what “type” your data is in R is fairly foreign to most R users, and misusing the **type** argument in **Rmpi** is a very easy way to crash your program. Even if you are more comfortable with statically typed languages and have no problem with this concept, consider the following example:

Types in R

```
1 > is.integer(1)
2 [1] FALSE
3 > is.integer(2)
4 [1] FALSE
5 > is.integer(1:2)
6 [1] TRUE
```

There are good reasons for R Core to have made this choice; that is not the point. The point is that because objects in R are dynamically typed, having to know the type of your data when utilizing **Rmpi** is a needless burden. Instead, **pbdMPI** takes the approach of adding a small abstraction layer on top (which we intend to show does not negatively impact performance) so that the user need not worry about such fiddly details.

In terms of performance, **pbdMPI** can greatly outperform **Rmpi**. We present here the results of a benchmark we performed comparing the “allgather” operation between the two packages (Schmidt et al., 2012e). The benchmark consisted of calling the respective “allgather” function from each package on a randomly generated $10,000 \times 10,000$ distributed matrix with entries coming from the standard normal distribution, using different numbers of processors. Table 3.1 shows the results for this test, and in each case, **pbdMPI** is the clear victor.

Whichever package you choose, whichever your favorite, for the remainder of this document we will be using (either implicitly or explicitly) **pbdMPI**.

¹See https://en.wikipedia.org/wiki/Code_golf

Table 3.1: Benchmark Comparing **Rmpi** and **pbmMPI**. Run time in seconds is listed for each operation. The speedup is relative to **Rmpi**.

| Cores | Rmpi | pbmMPI | Speedup |
|-------|-------------|---------------|---------|
| 32 | 24.6 | 6.7 | 3.67 |
| 64 | 25.2 | 7.1 | 3.55 |
| 128 | 22.3 | 7.2 | 3.10 |
| 256 | 22.4 | 7.1 | 3.15 |

3.3 The GBD Data Structure

This is the boring stuff everyone hates, but like your medicine, it's ultimately better for you to just take it and get it out of the way: data structures. In particular, we will be discussing a distributed data structure that, for lack of a better name (and I assure you are tried), we will call the GBD data structure. This data structure is more paradigm or philosophy than a rigid data structure like an array or list. Consider it a set of “best practices”, or if nothing else, a starting place if you have no idea how to proceed.

The GBD data structure is designed to fit the types of problems which are arguably most common to data science, namely tall and skinny matrices. It will work best with these (from a computational efficiency perspective) problems, although that is not required. In fact, very little at all is required of this data structure. At its core, the data structure is a distributed matrix data structure, with the following rules:

1. GBD is *distributed*. No one processor owns all of the matrix.
2. GBD is *non-overlapping*. Any row owned by one processor is owned by no other processors.
3. GBD is *row-contiguous*. If a processor owns one element of a row, it owns the entire row.
4. GBD is globally *row-major*², locally *column-major*³.
5. The last row of the local storage of a processor is adjacent (by global row) to the first row of the local storage of next processor (by communicator number) that owns data. That is, global row-adjacency is preserved in local storage via the communicator.
6. GBD is (relatively) easy to understand, but can lead to bottlenecks if you have many more columns than rows.

Of this list, perhaps the most difficult to understand is number 5. This is a precise, albeit cumbersome explanation for a simple idea. If two processors are adjacent and each owns data, then their local sub-matrices are adjacent row-wise as well. For example, rows n and $n + 1$ of a matrix are adjacent; possible configurations for the distributed ownership are processors q owns row n and $q + 1$ owns row $n + 1$; processor q owns row n , processor $q + 1$ owns *nothing*, and processor $q + 2$ owns row $n + 1$.

²In the sense of the data decomposition. More specifically, the global matrix is chopped up into local sub-matrices in a row-major way.

³The local sub-objects are R matrices, which are stored in column-major fashion.

For some, no matter how much we try or what we write, the wall of text simply will not suffice. So here are a few visual examples. Suppose we have the global data matrix x , given as:

$$x = \begin{bmatrix} x_{11} & x_{12} & x_{13} & x_{14} & x_{15} & x_{16} & x_{17} & x_{18} & x_{19} \\ x_{21} & x_{22} & x_{23} & x_{24} & x_{25} & x_{26} & x_{27} & x_{28} & x_{29} \\ x_{31} & x_{32} & x_{33} & x_{34} & x_{35} & x_{36} & x_{37} & x_{38} & x_{39} \\ x_{41} & x_{42} & x_{43} & x_{44} & x_{45} & x_{46} & x_{47} & x_{48} & x_{49} \\ x_{51} & x_{52} & x_{53} & x_{54} & x_{55} & x_{56} & x_{57} & x_{58} & x_{59} \\ x_{61} & x_{62} & x_{63} & x_{64} & x_{65} & x_{66} & x_{67} & x_{68} & x_{69} \\ x_{71} & x_{72} & x_{73} & x_{74} & x_{75} & x_{76} & x_{77} & x_{78} & x_{79} \\ x_{81} & x_{82} & x_{83} & x_{84} & x_{85} & x_{86} & x_{87} & x_{88} & x_{89} \\ x_{91} & x_{92} & x_{93} & x_{94} & x_{95} & x_{96} & x_{97} & x_{98} & x_{99} \end{bmatrix}_{9 \times 9}$$

with processor array⁴ (indexing always starts at 0 not 1)

$$\text{Processors} = \begin{matrix} 0 & 1 & 2 & 3 & 4 & 5 \end{matrix}$$

Then we might split up and distribute the data onto processors like so:

$$x = \begin{bmatrix} x_{11} & x_{12} & x_{13} & x_{14} & x_{15} & x_{16} & x_{17} & x_{18} & x_{19} \\ x_{21} & x_{22} & x_{23} & x_{24} & x_{25} & x_{26} & x_{27} & x_{28} & x_{29} \\ \hline x_{31} & x_{32} & x_{33} & x_{34} & x_{35} & x_{36} & x_{37} & x_{38} & x_{39} \\ x_{41} & x_{42} & x_{43} & x_{44} & x_{45} & x_{46} & x_{47} & x_{48} & x_{49} \\ \hline x_{51} & x_{52} & x_{53} & x_{54} & x_{55} & x_{56} & x_{57} & x_{58} & x_{59} \\ x_{61} & x_{62} & x_{63} & x_{64} & x_{65} & x_{66} & x_{67} & x_{68} & x_{69} \\ \hline x_{71} & x_{72} & x_{73} & x_{74} & x_{75} & x_{76} & x_{77} & x_{78} & x_{79} \\ x_{81} & x_{82} & x_{83} & x_{84} & x_{85} & x_{86} & x_{87} & x_{88} & x_{89} \\ \hline x_{91} & x_{92} & x_{93} & x_{94} & x_{95} & x_{96} & x_{97} & x_{98} & x_{99} \end{bmatrix}_{9 \times 9}$$

With local storage view:

$$\begin{bmatrix} x_{11} & x_{12} & x_{13} & x_{14} & x_{15} & x_{16} & x_{17} & x_{18} & x_{19} \\ x_{21} & x_{22} & x_{23} & x_{24} & x_{25} & x_{26} & x_{27} & x_{28} & x_{29} \end{bmatrix}_{2 \times 9}$$

$$\begin{bmatrix} x_{31} & x_{32} & x_{33} & x_{34} & x_{35} & x_{36} & x_{37} & x_{38} & x_{39} \\ x_{41} & x_{42} & x_{43} & x_{44} & x_{45} & x_{46} & x_{47} & x_{48} & x_{49} \end{bmatrix}_{2 \times 9}$$

$$\begin{bmatrix} x_{51} & x_{52} & x_{53} & x_{54} & x_{55} & x_{56} & x_{57} & x_{58} & x_{59} \\ x_{61} & x_{62} & x_{63} & x_{64} & x_{65} & x_{66} & x_{67} & x_{68} & x_{69} \end{bmatrix}_{2 \times 9}$$

$$\begin{bmatrix} x_{71} & x_{72} & x_{73} & x_{74} & x_{75} & x_{76} & x_{77} & x_{78} & x_{79} \end{bmatrix}_{1 \times 9}$$

$$\begin{bmatrix} x_{81} & x_{82} & x_{83} & x_{84} & x_{85} & x_{86} & x_{87} & x_{88} & x_{89} \end{bmatrix}_{1 \times 9}$$

$$\begin{bmatrix} x_{91} & x_{92} & x_{93} & x_{94} & x_{95} & x_{96} & x_{97} & x_{98} & x_{99} \end{bmatrix}_{1 \times 9}$$

⁴Palette selected to be distinguishable by the color blind, taken from <http://jfly.iam.u-tokyo.ac.jp/color/#pallet>

This is a *load balanced* approach, where we try to give each processor roughly the same amount of stuff. Of course, that is not part of the rules of the GBD structure, so we could just as well distribute the data like so:

$$x = \left[\begin{array}{ccccccccc} x_{11} & x_{12} & x_{13} & x_{14} & x_{15} & x_{16} & x_{17} & x_{18} & x_{19} \\ x_{21} & x_{22} & x_{23} & x_{24} & x_{25} & x_{26} & x_{27} & x_{28} & x_{29} \\ x_{31} & x_{32} & x_{33} & x_{34} & x_{35} & x_{36} & x_{37} & x_{38} & x_{39} \\ x_{41} & x_{42} & x_{43} & x_{44} & x_{45} & x_{46} & x_{47} & x_{48} & x_{49} \\ \hline x_{51} & x_{52} & x_{53} & x_{54} & x_{55} & x_{56} & x_{57} & x_{58} & x_{59} \\ x_{61} & x_{62} & x_{63} & x_{64} & x_{65} & x_{66} & x_{67} & x_{68} & x_{69} \\ \hline x_{71} & x_{72} & x_{73} & x_{74} & x_{75} & x_{76} & x_{77} & x_{78} & x_{79} \\ \hline & & & & & & & & \\ \hline x_{81} & x_{82} & x_{83} & x_{84} & x_{85} & x_{86} & x_{87} & x_{88} & x_{89} \\ x_{91} & x_{92} & x_{93} & x_{94} & x_{95} & x_{96} & x_{97} & x_{98} & x_{99} \end{array} \right]_{9 \times 9}$$

With local storage view:

$$\left[\begin{array}{ccccccccc} & & & & & & & & \\ & & & & & & & & \\ & & & & & & & & \\ & & & & & & & & \\ & & & & & & & & \\ & & & & & & & & \\ & & & & & & & & \\ & & & & & & & & \\ & & & & & & & & \end{array} \right]_{0 \times 9}$$

$$\left[\begin{array}{ccccccccc} x_{11} & x_{12} & x_{13} & x_{14} & x_{15} & x_{16} & x_{17} & x_{18} & x_{19} \\ x_{21} & x_{22} & x_{23} & x_{24} & x_{25} & x_{26} & x_{27} & x_{28} & x_{29} \\ x_{31} & x_{32} & x_{33} & x_{34} & x_{35} & x_{36} & x_{37} & x_{38} & x_{39} \\ x_{41} & x_{42} & x_{43} & x_{44} & x_{45} & x_{46} & x_{47} & x_{48} & x_{49} \end{array} \right]_{4 \times 9}$$

$$\left[\begin{array}{ccccccccc} x_{51} & x_{52} & x_{53} & x_{54} & x_{55} & x_{56} & x_{57} & x_{58} & x_{59} \\ x_{61} & x_{62} & x_{63} & x_{64} & x_{65} & x_{66} & x_{67} & x_{68} & x_{69} \end{array} \right]_{2 \times 9}$$

$$\left[\begin{array}{ccccccccc} x_{71} & x_{72} & x_{73} & x_{74} & x_{75} & x_{76} & x_{77} & x_{78} & x_{79} \end{array} \right]_{1 \times 9}$$

$$\left[\begin{array}{ccccccccc} & & & & & & & & \\ & & & & & & & & \\ & & & & & & & & \\ & & & & & & & & \\ & & & & & & & & \end{array} \right]_{0 \times 9}$$

$$\left[\begin{array}{ccccccccc} x_{81} & x_{82} & x_{83} & x_{84} & x_{85} & x_{86} & x_{87} & x_{88} & x_{89} \\ x_{91} & x_{92} & x_{93} & x_{94} & x_{95} & x_{96} & x_{97} & x_{98} & x_{99} \end{array} \right]_{2 \times 9}$$

Generally, we would recommend using a load balanced approach over this bizarre distribution, although some problems may call for very strange data distributions. For example, it is possible and common to have an empty matrix after some subsetting or selection.

With our first of two cumbersome data structures out of the way, we can proceed to much more interesting content: actually using MPI.

3.4 Common MPI Operations

Fully explaining the process of MPI programming is a daunting task. Thankfully, we can punt and merely highlight some key MPI operations and how one should use them with **pbdMPI**.

3.4.1 Basic Communicator Wrangling

First things first, we must examine basic communicator issues, like construction, destruction, and each processor's position within a communicator.

- **Managing a Communicator:** Create and destroy communicators.
`init()` — initialize communicator
`finalize()` — shut down communicator(s)
- **Rank query:** Determine the processor's position in the communicator.
`comm.rank()` — “who am I?”
`comm.size()` — “how many of us are there?”
- **Barrier:** No processor can proceed until *all* processors can proceed.
`barrier()` — “computation wall” that only all processors together can tear down.

One quick word before proceeding. If a processor queries `comm.size()`, this will return the total number of processors in the communicators. However, communicator indexing is like indexing in the programming language C. That is, the first element is numbered 0 rather than 1. So when the first processor queries `comm.rank()`, it will return 0, and when the last processor queries `comm.rank()`, it will return `comm.size() - 1`.

We are finally ready to write our first MPI program:

Simple pbdMPI Example 1

```
1 library(pbdMPI, quiet = TRUE)
2 init()
3
4 myRank <- comm.rank() + 1 # comm index starts at 0, not 1
5 print(myRank)
6
7 finalize()
```

Unfortunately, it is not very exciting, but you have to crawl before you can drag race. Remember that all of our programs are written in SPMD style. So this *one* single program is written, and each processor will execute the same program, but with different results, whence the name “Single Program/Multiple Data”.

So what does it do? First we initialize the MPI communicator with the call to `init()`. Next, we have each processor query its rank via `comm.rank()`. Since indexing of MPI communicators starts at 0, we add 1 because that is what we felt like doing. Finally we call R's `print()` function to print the result. This printing is not particularly clever, and each processor will be clamoring to dump its result to the output file/terminal. We will discuss more sophisticated means of printing later. Finally, we shut down the MPI communicator with `finalize()`.

If you were to save this program in the file `mpiex1.r` and you wished to run it with 2 processors, you would issue the command:

Shell Command

```
### (Use Rscript.exe for windows system)
mpiexec -np 2 Rscript mpiex1.r
```

To use more processors, you modify the `-np` argument (“number processors”). So if you want to use 4, you pass `-np 4`.

The above program technically, though not in spirit, bucks the trend of officially opening with a “Hello World” program. So as not to incur the wrath of the programming gods, we offer a simple such example by slightly modifying the above program:

Simple pbdMPI Example 1.5

```
1 library(pbdMPI, quiet = TRUE)
2 init()
3
4 myRank <- comm.rank()
5
6 if (myRank == 0){
7   print("Hello, world.")
8 }
9
10 finalize()
```

One word of general warning we offer now is that we are taking very simple approaches here for the sake of understanding, heavily relying on function argument defaults. However, there are all kinds of crazy, needlessly complicated things you can do with these functions. See the **pbdMPI** reference manual for full details about how one may utilize these (and other) **pbdMPI** functions.

3.4.2 Reduce, Broadcast, and Gather

Once managing a communicator is under control, you presumably want to do things with all of your processors. The typical way you will have the processors interact is given below:

- **Reduction:** Say each processor has a number `x.gbd`. Add all of them up, find the largest, find the smallest,
`reduce(x.gbd, op='sum')` — only one processor gets result (default is 0)
`allreduce(x.gbd, op='sum')` — every processor gets result
- **Gather:** Say each processor has a number. Create a new object on some processor(s) containing all of those numbers.
`gather(x.gbd)` — only one processor gets result
`allgather(x.gbd)` — every processor gets result
- **Broadcast:** One processor has a number `x.gbd` that every other processor should also have.
`bcast(x.gbd)`

Here perhaps explanations are inferior to examples; so without wasting any more time, we proceed to the next example:

Simple pbdMPI Example 2

```
1 library(pbdMPI, quiet = TRUE)
2 init()
3
4 n.gbd <- sample(1:10, size=1)
5
6 sm <- allreduce(n.gbd) # default op is 'sum'
7 print(sm)
8
9 gt <- allgather(n.gbd)
10 print(gt)
11
12 finalize()
```

So what does it do? First each processor samples a number from 1 to 10; it is probably true that each processor will be using a different seed for this sampling, though you should not rely on this alone to ensure good parallel seeds. More on this particular problem in Section 3.4.3 below.

Next, we perform an `allreduce()` on `n.gbd`. Conceivably, the processors could have different values for `n.gbd`. So the value of `n` is local to each processor. So perhaps we want to add up all these numbers (with as many numbers as there are processors) and store them in the global value `sm` (for “sum”). Each processor will agree as to the value of `sm`, even if they disagree about the value of `n.gbd`.

Finally, we do the same but with an `allgather()` operation.

Try experimenting with this by running the program several times. You should get different results each time. To make sure we have not been lying to you about what is happening, you can even print the values of `n.gbd` before the reduce and gather operations.

3.4.3 Printing and RNG Seeds

In addition to the above common MPI operations, which will make up the bulk of the MPI programmer’s toolbox, we offer a few extra utility functions:

- **Print:** printing with control over which processor prints.
`comm.print(x, ...)`
`comm.cat(x, ...)`
- **Random Seeds:**
`comm.set.seed(seed, diff=FALSE)`: every processor uses the seed `seed`
`comm.set.seed(seed, diff=TRUE)`: every processor uses an independent seed (via **rlcuyer**)

The `comm.print()` and `comm.cat()` functions are especially handy, because they are much more sophisticated than their R counterparts when using multiple processes. These functions which processes do the printing, and if you choose to have all processes print their result, then the printing occurs in an orderly fashion, with processor 0 getting the first line, processor 1 getting the second, and so on.

For example, revisiting our “Hello, world” example, we can somewhat simplify the code with a slight modification:

Simple pbdMPI Example 3

```
1 library(pbdMPI, quiet = TRUE)
2 init()
3
4 myRank <- comm.rank()
5
6 comm.print("Hello, world.")
7
8 finalize()
```

If we want to see what each processor has to say, we can pass the optional argument `all,rank=TRUE` to `comm.print()`. By default, each process will print its rank, then what you told it to print. You can suppress the printing of communicator rank via the optional argument `quiet=TRUE` to `comm.print()`.

These functions are quite handy...

HOWEVER

these functions are potentially dangerous, and so some degree of care should be exercised. Indeed, it is possible to lock up all of the active R sessions by incorrectly using them. Worse, achieving this behavior is fairly easy to do. The way this occurs is by issuing a `comm.print()` on an expression which requires communication. For example, suppose we have a distributed object with local piece `x.gbd` and a function `myFunction()` which requires communication between the processors. Then calling

A Cautionary Tale of Printing in Parallel (1 of 3)

```
1 print(myFunction(x.gbd))
```

is just fine, but will not have the nice orderly, behaved printing style of `comm.print()`. However, if we issue

A Cautionary Tale of Printing in Parallel (2 of 3)

```
1 comm.print(myFunction(x.gbd))
```

then we have just locked up all of the R processes. Indeed, behind the scenes, a call somewhat akin to

```

1 for (rank in 0:comm.size()){
2   if (comm.rank() == rank){
3     # do things
4   }
5   barrier()
6 }

```

has been ordered. The problem arises in the “do things” part. Since (in our hypothetical example) the function `myFunction()` requires communication between the processors, it will simply wait forever for the others to respond until the job is killed. This is because the other processors skipped over the “do things” part and are waiting at the barrier. So lonely little processor 0 has been stood up, unable to communicate with the remaining processors.

To avoid this problem, make it a personal habit to only print on *results*, not *computations*. We can quickly rectify the above example by doing the following:

A Cautionary Tale of Printing in Parallel (3 of 3)

```

1 myResult <- myFunction(x.gbd)
2 comm.print(myResult)

```

In short, printing stored objects is safe. Printing a yet-to-be-evaluated expression is not safe.

3.4.4 Apply, Lapply, and Sapply

But the **pbdMPI** sugar extends to more than just printing. We also have a family of “*ply” functions, in the same vein as R’s `apply()`, `lapply()`, and `sapply()`:

- **Apply:** *ply-like functions.
 - `pbdApply(X, MARGIN, FUN, ...)` — analogue of `apply()`
 - `pbdLapply(X, FUN, ...)` — analogue of `lapply()`
 - `pbdSapply(X, FUN, ...)` — analogue of `sapply()`

For more efficient approach (non-barrier), one may consider use task pull parallelism instead of “*ply” functions, see Section 14.4 for more details.

Here is a simple example utilizing `pbdLapply()`:

Example 4

```

1 library(pbdMPI, quiet = TRUE)
2 init()
3
4 n <- 100
5 x <- split((1:n) + n * comm.rank(), rep(1:10, each = 10))
6 sm <- pbdLapply(x, sum)
7 comm.print(unlist(sm))

```



```
8  
9 finalize()
```

So what does it do? Why don't you tell us? We're busy people, after all, and we're not going to be around forever. Try guessing what it will do, then run the program to see if you are correct. As you evaluate this and every parallel code, ask yourself which pieces involve communication and which pieces are local computations.

3.5 Miscellaneous Basic MPI Tasks

3.5.1 Timing MPI Tasks

Measuring run time is a fundamental performance measure in computing. However, in parallel computing, not all “parallel components” (e.g. threads, or MPI processes) will take the same amount of time to complete a task, even when all tasks are given completely identical jobs. So measuring “total run time” begs the question, run time of what?

To help, we offer a timing function `timer()` which can wrap segments of code much in the same way that `system.time()` does. However, the three numbers reported by `timer()` are:

- the minimum elapsed time measured across all processes,
- the average elapsed time measured across all processes, and
- the maximum elapsed time across all processes.

The code for this function is listed below:

Timer Function

```
1 timer <- function(timed)  
2 {  
3   ltime <- system.time(timed)[3]  
4  
5   mintime <- allreduce(ltime, op='min')  
6   maxtime <- allreduce(ltime, op='max')  
7  
8   meantime <- allreduce(ltime, op='sum') / comm.size()  
9  
10  return( c(min=mintime, mean=meantime, max=maxtime) )  
11 }
```

3.5.2 Distributed Logic

Example: Manage comparisons across all MPI processes.

The demo command is

```
### At the shell prompt, run the demo with 4 processors by
### (Use Rscript.exe for windows system)
mpiexec -np 4 Rscript -e "demo(comparators,'pbdDEMO',ask=F,echo=F)"
```

This final MPI example is not statistical in nature, but is very useful all the same, and so we include it here. The case frequently arises where the MPI programmer will need to do logical comparisons across all processes. The idea is to extend the very handy `all()` and `any()` base R functions to operate similarly on distributed logicals.

You could do this directly. Say you want to see if any processes have `TRUE` stored in the variable `localLogical`. This amounts to something on the order of:

R Code

```
1 globalLogical <- as.logical(allreduce(localLogical, op='max'))
```

Or you can use the function `comm.any()` from **pbdMPI**:

R Code

```
1 globalLogical <- comm.any(localLogical)
```

which essentially does the same thing, but is more concise. Likewise, there is a `comm.all()` function, which in the equivalent “long-form” above would use `op='min'`.

The demo for these functions consists of two parts. For the first, we do a simple demonstration of how these functions behave:

R Code

```
1 rank <- comm.rank()
2
3 comm.cat("\ntest value:\n", quiet=T)
4 test <- (rank > 0)
5 comm.print(test, all.rank=T, quiet=T)
6
7 comm.cat("\ncomm.all:\n", quiet=T)
8 test.all <- comm.all(test)
9 comm.print(test.all, all.rank=T, quiet=T)
10
11 comm.cat("\ncomm.any:\n", quiet=T)
12 test.any <- comm.any(test)
13 comm.print(test.any, all.rank=T, quiet=T)
```

which should have the output:

```
test value:
[1] FALSE
[1] TRUE
[1] TRUE
```

```
[1] TRUE

comm.all:
[1] FALSE
[1] FALSE
[1] FALSE
[1] FALSE

comm.any:
[1] TRUE
[1] TRUE
[1] TRUE
[1] TRUE
```

The demo also has another use case which could be very useful to a developer. You may be interested in trying something on only one processor and then shutting down all MPI processes if problems are encountered. To do this in SPMD style, you can create a variable on all processes to track whether a problem has been encountered. Then after critical code sections, use `comm.any()` to update and act appropriately. A very simple example is provided below.

R Code

```
1 need2stop <- FALSE
2
3 if (rank==0){
4   need2stop <- TRUE
5 }
6
7 need2stop <- comm.any(need2stop)
8
9 if (need2stop)
10  stop("Problem :[")
```

3.6 Exercises

- 3-1 Write a script that will have each processor randomly take a sample of size 1 of `TRUE` and `FALSE`. Have each processor print its result.
- 3-2 Modify the script in Exercise 3-1 above to determine if any processors sampled `TRUE`. Do the same to determine if all processors sampled `TRUE`. In each case, print the result. Compare to the functions `comm.all()` and `comm.any()`. [Hint: use `allreduce\(\)`](#).
- 3-3 In **pbdMPI**, there is a parallel sorting function called `comm.sort()` which is similar to the usual `sort()` of R. Implement parallel equivalents to the usual `order()` and `rank()` of R.
- 3-4 Time the performance of Exercise 3-3. Identify the need of MPI communications for different sorting or ordering algorithms.

- 3-5 There are “parallel copycat” versions of to R’s *ply functions in several R packages, such as `mclapply()` (a parallel version of `lapply()`) in the **parallel** package. Try to compare the difference and performance of those *ply-like functions.
- 3-6 Be aware that in SPMD programming, calling stack in-balance may result crashes, message truncations, and wrong results especially when condition operations involved with MPI communications. Suggesting correct ways to the following example.

R Code

```

1 if(comm.rank() == 0){
2   ret <- reduce(1)
3   ret.1 <- allreduce(2)
4 } else{
5   ret <- allreduce(2)
6   ret.1 <- reduce(1)
7 }

```

- 3-7 Also, in SPMD programming, point to point communications are the basic way to send and receive data from other processors. However, blocked calls may result dead locks of entire MPI communicator, then hang R programs without further warnings. Suggesting correct ways to the following example.

R Code

```

1 if(comm.rank() == 0){
2   ret <- send(1, rank.dest = 1)
3 } else if(comm.rank() == 1){
4   ret <- recv(2, rank.source = 0)
5 }

```

- 3-8 Further, in SPMD programming, data type and data size to MPI communicating functions are very sensitive especially when lower level functions involved for gaining performance. On other hand, R has less sensitive to data type and data size, and sometimes it converts data internally. Suggesting correct ways to the following example.

R Code

```

1 if(comm.rank() == 0){
2   ret <- spmd.send.double(1:2)
3 } else{
4   ret <- spmd.recv.double(2.0)
5 }

```

- 3-9 Note that unprotected objects in R may be recycled (garbage collection, see `gc()` for details) at any time from memory. This is a nice mechanism for memory managements; however, it also issues a problem for non-blocked communication which can avoid dead lock problems. The problem typically occurs when sending large number of small buffers where unprotected buffers can either be released from R or overwrite by other R objects before actual communications are finished. Suggesting better ways to the following example.

R Code

```
1  if(comm.rank() == 0){  
2    for(i in 1:100){  
3      ret <- spmd.isend.integer(i, rank.dest = 1)  
4    }  
5    ### Further computation unreleated to "i".  
6    wait()  
7  } else if(comm.rank() == 1){  
8    for(i in 1:100){  
9      ret <- spmd.recv.integer(i)  
10     print(ret)  
11   }  
12 }
```

Basic Statistics Examples

And perhaps, posterity will thank me for having shown it that the ancients did not know everything.

—Pierre de Fermat

This chapter introduces a few simple examples and explains a little about computing with distributed data directly over MPI. These implemented examples/functions are partly selected from the Cookbook of HPSC website (Chen and Ostrouchov, 2011) at <http://thirteen-01.stat.iastate.edu/snoweye/hpsc/?item=cookbook>.

4.1 Monte Carlo Simulation

Example: Compute a numerical approximation for π .

The demo command is

```
### At the shell prompt, run the demo with 4 processors by
### (Use Rscript.exe for windows system)
mpiexec -np 4 Rscript -e "demo(monte_carlo,'pbdDEMO',ask=F,echo=F)"
```

This is a simple Monte Carlo simulation example for numerically estimating π . Suppose we sample N uniform observations (x_i, y_i) inside (or perhaps on the border of) the unit square $[0, 1] \times [0, 1]$, where $i = 1, 2, \dots, N$. Then

$$\pi \approx 4 \frac{L}{N} \tag{4.1}$$

where $0 \leq L \leq N$ is the number of observations sampled satisfying

$$x_i^2 + y_i^2 \leq 1 \tag{4.2}$$

The intuitive explanation for this strategy which is sometimes given belies a misunderstanding of infinite cardinalities, and infinite processes in general. We are not *directly* approximating an

area through this sampling scheme, because to do so with a finite-point sampling scheme would be madness requiring a transfinite process. Indeed, let S_N be the collection of elements satisfying inequality (4.2). Then note that for each $N \in \mathbb{N}$ that the area of S_N is precisely 0. Whence,

$$\lim_{N \rightarrow \infty} \text{Area}(S_N) = 0$$

This bears repeating. Finite sampling of an uncountable space requires uncountably many such sampling operations to “fill” the infinite space. For a proper treatment of set theoretic constructions, including infinite cardinals, see (Kunen, 1980).

One could argue that we are evaluating a ratio of integrals with each using the counting measure, which satisfies technical correctness but is far from clear. Now indeed, certain facts of area are vital here, but some care should be taken in the discussion as to what exactly justifies our claim in (4.1).

In reality, we are evaluating the probability that someone throwing a 0-dimensional “dart” at the unit square will have that “dart” also land below the arc of the unit circle contained within the unit square. Formally, let U_1 and U_2 be random uniform variables, each from the closed unit interval $[0, 1]$. Define the random variable

$$X := \begin{cases} 1, & U_1^2 + U_2^2 \leq 1 \\ 0, & \text{otherwise} \end{cases}$$

Let $V_i = U_i^2$ for $i = 1, 2$. Then the expected value

$$\begin{aligned} \mathbb{E}[X] &= \mathbb{P}(V_1 + V_2 \leq 1) \\ &= \int_0^1 \int_0^{1-V_1} p(V_1, V_2) dV_2 dV_1 \\ &= \int_0^1 \int_0^{1-V_1} \left(\frac{1}{2\sqrt{V_1}} \right) \left(\frac{1}{2\sqrt{V_2}} \right) dV_2 dV_1 \\ &= \frac{1}{2} \int_0^1 \left(\frac{1-V_1}{V_1} \right)^{1/2} dV_1 \\ &= \frac{1}{2} \left[V_1 \left(\frac{1-V_1}{V_1} \right)^{1/2} - \frac{1}{2} \arctan \left(\frac{\left(\frac{1-V_1}{V_1} \right)^{1/2} (2V_1 - 1)}{2(V_1 - 1)} \right) \right]_{V_1 \rightarrow 0}^{V_1 \rightarrow 1} \\ &= \frac{1}{2} \left[\frac{\pi}{4} + \frac{\pi}{4} \right] \end{aligned}$$

and by sampling observations X_i for $i = 1, \dots, N$, by the Strong Law of Large Numbers

$$\bar{X}_N \xrightarrow{a.s.} \frac{\pi}{4} \quad \text{as } N \rightarrow \infty \quad (4.3)$$

In other words,

$$\mathbb{P} \left(\lim_{N \rightarrow \infty} \bar{X}_N = \frac{\pi}{4} \right) = 1$$

Whence,

$$\frac{L}{N} \xrightarrow{a.s.} \frac{\pi}{4} \quad \text{as } N \rightarrow \infty$$

But because no one is going to read that, and if they do they'll just call the author a grumpy old man, the misleading picture you desire can be found in Figure 4.1. And to everyone who

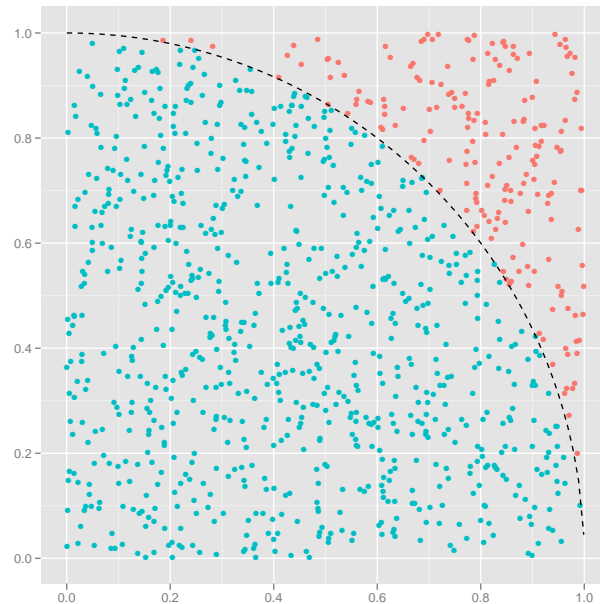


Figure 4.1: Approximating π by Monte Carlo methods

found this looking for a homework solution, you're welcome.

The key step of the demo code is in the following block:

R Code

```

1 N.gbd <- 1000
2 X.gbd <- matrix(runif(N.gbd * 2), ncol = 2)
3 r.gbd <- sum(rowSums(X.gbd^2) <= 1)
4 ret <- allreduce(c(N.gbd, r.gbd), op = "sum")
5 PI <- 4 * ret[2] / ret[1]
6 comm.print(PI)

```

In line 1, we specify sample size in `N.gbd` for each processor, and $N = D \times \text{N.gbd}$ if D processors are executed. In line 2, we generate samples in `X.gbd` for every processor. In line 3, we compute how many of the “radii” are less than or equal to 1 for each processors. In line 4, we call `allreduce()` to obtain total numbers across all processors. In line 5, we use the Equation (4.1). Since SPMD, `ret` is common on all processors, and so is `PI`.

4.2 Sample Mean and Sample Variance

Example: Compute sample mean/variance for distributed data.

The demo command is

```
### At the shell prompt, run the demo with 4 processors by
### (Use Rscript.exe for windows system)
mpiexec -np 4 Rscript -e "demo(sample_stat,'pbdDEMO',ask=F,echo=F)"
```

Suppose $\mathbf{x} = \{x_1, x_2, \dots, x_N\}$ are observed samples, and N is potentially very large. We can distribute \mathbf{x} in 4 processors, and each processor receives a proportional amount of data. One simple way to compute sample mean \bar{x} and sample variance s_x is based on the formulas:

$$\begin{aligned}\bar{x} &= \frac{1}{N} \sum_{n=1}^N x_n \\ &= \sum_{n=1}^N \frac{x_n}{N}\end{aligned}\tag{4.4}$$

and

$$\begin{aligned}s_x &= \frac{1}{N-1} \sum_{n=1}^N (x_n - \bar{x})^2 \\ &= \frac{1}{N-1} \sum_{n=1}^N x_n^2 - \frac{2\bar{x}}{N-1} \sum_{n=1}^N x_n + \frac{1}{N-1} \sum_{n=1}^N \bar{x}^2 \\ &= \sum_{n=1}^N \left(\frac{x_n^2}{N-1} \right) - \frac{N\bar{x}^2}{N-1}\end{aligned}\tag{4.5}$$

where expressions (4.4) and (4.5) are one-pass algorithms, which are potentially faster than the first expressions, especially for large N . However, this method of computing the variance in one pass can suffer from round-off errors, and so in general is not numerically stable. We provide this here for demonstration purposes only. Additionally, only the first and second moments are implemented, while the extension of one-pass algorithms to higher order moments is also possible.

The demo generates random data on 4 processors, then utilizes the `mpi.stat()` function:

R Code

```
1 mpi.stat <- function(x.gbd){
2   ### For mean(x).
3   N <- allreduce(length(x.gbd), op = "sum")
4   bar.x.gbd <- sum(x.gbd / N)
```

```

5  bar.x <- allreduce(bar.x.gbd, op = "sum")
6
7  ### For var(x).
8  s.x.gbd <- sum(x.gbd^2 / (N - 1))
9  s.x <- allreduce(s.x.gbd, op = "sum") - bar.x^2 * (N / (N - 1))
10
11  list(mean = bar.x, s = s.x)
12 } # End of mpi.stat().

```

where `allreduce()` in **pbdMPI** (Chen et al., 2012a) can be utilized in this examples to aggregate local information across all processors.

4.3 Binning

Example: Find binning counts for distributed data.

The demo command is

```

### At the shell prompt, run the demo with 4 processors by
### (Use Rscript.exe for windows system)
mpirun -np 4 Rscript -e "demo(binning, 'pbdDEMO', ask=F, echo=F)"

```

Binning is a classical problem in statistics which helps to quickly summarize the data structure by setting some “breaks” between the minimum and maximum values. This is a particularly useful tool for constructing histograms, as well as categorical data analysis.

The demo generates random data on 4 processors, then utilizes the `mpi.bin()` function:

R Code

```

1  mpi.bin <- function(x.gbd, breaks = pi / 3 * (-3:3)){
2    bin.gbd <- table(cut(x.gbd, breaks = breaks))
3    bin <- as.array(allreduce(bin.gbd, op = "sum"))
4    dimnames(bin) <- dimnames(bin.gbd)
5    class(bin) <- class(bin.gbd)
6    bin
7  } # End of mpi.bin().

```

This simple implementation utilizes R’s own `table()` function to obtain local counts, then calls `allreduce()` to obtain global counts on all processors.

4.4 Quantile

Example: Compute sample quantile order statistics for distributed data.

The demo command is

```
### At the shell prompt, run the demo with 4 processors by
### (Use Rscript.exe for windows system)
mpiexec -np 4 Rscript -e "demo(quantile,'pbdDEMO',ask=F,echo=F)"
```

Another fundamental tool in the statistician's toolbox is finding quantiles. Quantiles are points taken from the cumulative distribution function. Formally, a q -quantile (or q -tile) with $q \in [0, 1]$ of a random variable X is any value θ_q such that¹

$$\begin{aligned} \mathbb{P}(X \leq \theta_q) &\geq q & \text{and} \\ \mathbb{P}(X \geq \theta_q) &\leq 1 - q \end{aligned}$$

Note that by this definition, a quantile neither need exist or be unique. Indeed, for the former, consider the standard normal distribution with $q = 1$, and for the latter consider the probability 0 values of a uniform distribution. Perhaps to narrow the scope of these problems, another common definition is

$$\theta_q = \inf\{x \mid \mathbb{P}(X \leq x) \geq q\}$$

In this example, we will be estimating quantiles from a sample. Doing so requires sub-dividing the data into q (almost) evenly sized subsets, giving rise to the language k 'th q -tile, for integers $0 < k < \frac{1}{q}$.

Before proceeding, we wish to make very clear the distinction between a theoretical quantile and quantile estimation, as many web pages confuse these two topics. A quantile estimation from a sample requires ordering and can take many forms; in fact, there are nine possible such forms in R's own `quantile()` function (see `help(quantile)` in R). The definitions of Kendall and Cramer may be the source of all the confusion (Benson, 1949). Kendall's definition, conflating the term "quantile" with the act of quantile estimation, seems to have entered most dictionaries (and Wikipedia), whereas mathematical statistics favors the more general and simple definition of Cramer.

This example can be extended to construct Q-Q plots, compute cumulative density function estimates and nonparametric statistics, as well as solve maximum likelihood estimators. This is perhaps an inefficient implementation to approximate a quantile and is not equivalent to the original `quantile()` function in R. But in some sense, it should work well at a large scale. The demo generates random data on 4 processors, then utilizes the `mpi.quantile()`:

R Code

```
1 mpi.quantile <- function(x.gbd, prob = 0.5){
2   if(sum(prob < 0 | prob > 1) > 0){
3     stop("prob should be in (0, 1)")
4   }
5
6   N <- allreduce(length(x.gbd), op = "sum")
7   x.max <- allreduce(max(x.gbd), op = "max")
```

¹This definition is due to the mathematical statistician Herman Rubin: <http://mathforum.org/kb/message.jspa?messageID=406278>

```

8  x.min <- allreduce(min(x.gbd), op = "min")
9
10 f.quantile <- function(x, prob = 0.5){
11   allreduce(sum(x.gbd <= x), op = "sum") / N - prob
12 }
13
14 uniroot(f.quantile, c(x.min, x.max), prob = prob[1])$root
15 } # End of mpi.quantile().

```

Here, a numerical function is solved by using `uniroot()` to find out the appropriate value where the cumulative probability is less than or equal to the specified quantile. Specifically, it finds the zero, or root, of the monotone `f.quantile()` function. This simple example shows that with just a little effort, direct MPI methods are greatly applicable on large scale data analysis and likelihood computing.

Note that in the way that the `uniroot()` call is used above, we are legitimately operating in parallel and on distributed data. Other optimization functions such as `optim()` and `nlm()` can be utilized in the same way.

4.5 Ordinary Least Squares

Example: Compute ordinary least square solutions for GBD distributed data.

The demo command is

```

### At the shell prompt, run the demo with 4 processors by
### (Use Rscript.exe for windows system)
mpirexec -np 4 Rscript -e "demo(ols,'pbdDEMO',ask=F,echo=F)"

```

Ordinary least squares (OLS) is perhaps *the* fundamental tool of the statistician. The goal is to find a solution β such that

$$\|X\beta - y\|_2^2 \quad (4.6)$$

is minimized. In statistics, we tend to prefer to think of the problem as being of the form

$$y = X\beta + \epsilon \quad (4.7)$$

where y is $N \times 1$ observed vector, X is $N \times p$ (possibly designed) matrix which is often assumed to have full rank (more on that later), and $N \gg p$, β is the unknown parameter to be estimated, and ϵ is errors and to be minimized in norm.

Note that above, we do indeed mean (in fact, stress) *a* solution to the linear least squares problem. For many applications a statistician will face, expression (4.6) will actually have a unique solution. But this is not always the case, and trouble often arises when the model matrix is rank-deficient. Indeed, in this case it may occur that there is an infinite family of solutions. So typically we go further and demand that a solution β be such that $\|\beta\|_2$ is at least as small as the corresponding norm of any other solution (although even this may not guarantee uniqueness).

A properly thorough treatment of the problems involved here go beyond the scope of this document, and require the reader have in-depth familiarity with linear algebra. For our purposes, the concise explanation above will suffice.

In the full rank case, we can provide an analytical, “closed-form” solution to the problem. In this case, the classical is given by:

$$\hat{\beta}_{ols} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad (4.8)$$

This example can be also generalized to weighted least squares (WLS), and linear mixed effect models. See http://en.wikipedia.org/wiki/Least_squares and http://en.wikipedia.org/wiki/Mixed_model for more details.

The implementation is straight forward:

R Code

```

1 if(length(y.gbd) != nrow(X.gbd)){
2   stop("length(y.gbd) != nrow(X.gbd)")
3 }
4
5 t.X.gbd <- t(X.gbd)
6 A <- allreduce(t.X.gbd %*% X.gbd, op = "sum")
7 B <- allreduce(t.X.gbd %*% y.gbd, op = "sum")
8
9 solve(matrix(A, ncol = ncol(X.gbd))) %*% B

```

While this is a fine demonstration of the power of “getting your hands dirty”, this approach is only efficient for small N and small p . This is, in large part, because the operation is not “fully parallel”, in that the solution is serial and replicated on all processors. Worse, directly computing

$$(\mathbf{X}^T \mathbf{X})^{-1}$$

has numerical stability issues. Instead, it is generally better (although much slower) to take an orthogonal factorization of the data matrix. See Appendix A for details.

Finally, all of the above assumes that the model matrix \mathbf{X} is full rank. However, we have implemented an efficient method of solving linear least squares problems in **pbdDMAT**’s `lm.fit()` method for distributed matrices. This method uses a fully parallel rank-revealing QR Decomposition to find the least squares solution. So for larger problems, and especially those where numerical accuracy is important or rank-degeneracy is a possibility, it is much better to simply convert `y.gbd` and `X.gbd` into the block-cyclic format as in Part III and utilize **pbdBASE** and **pbdDMAT** for all matrix computations.

4.6 Exercises

- 4-1 What are the assumptions used in order to invoke the Strong Law of Large Numbers (SLLN) in Statement (4.3)?

- 4-2 What is the Weak Law of Large Numbers (WLLN)? Prove that the SLLN implies the WLLN. Provide a counter example that the WLLN does not imply the SLLN.
- 4-3 In Statement (4.3), we showed that \bar{X}_N converges to $\frac{\pi}{4}$ almost surely, a very strong form of convergence. Show that additionally, \bar{X}_N converges to $\frac{\pi}{4}$ in probability by the WLLN, and that the sequence converges in distribution. (This can be as simple or as complicated as you like, depending on how many big theorems you wish to invoke).
- 4-4 Let $g : [0, 1] \rightarrow \mathbb{R}$ be a continuous function, and let \bar{X}_N be as in Statement (4.3). Show that $g(\bar{X}_N)$ converges to $g(\frac{\pi}{4})$ almost surely [Hint: use the property of continuity with respect to limits of sequences and the definition of almost sure convergence](#).
- 4-5 What are assumptions for Statement (4.7)? [Hint: Gauss-Markov Theorem](#).
- 4-6 Prove that $\hat{\beta}_{ols}$ of Statement (4.8) is an unbiased estimator of β provided appropriate assumptions, i.e., show that $\mathbb{E}[\hat{\beta}_{ols}] = \beta$.
- 4-7 Prove $\mathbf{X}^\top \mathbf{X}$ is non-negative definite if \mathbf{X} has full column rank p (and whence in this case, the inverse exists).
- 4-8 Iteratively Reweighted Least Squares (IRLS) is an important method for finding solutions to generalized linear models (GLM)². A common application of GLM's is logistic regression³. Implement a (not necessarily numerically stable) logistic regression function using IRLS for GBD data. For simplicity, you may wish to assume that the weighted matrix $\mathbf{X}^T \mathbf{W} \mathbf{X}$ is full rank at each iteration. [Hint: McCullagh and Nelder \(1989\)](#).

²See http://en.wikipedia.org/wiki/Generalized_linear_model for details

³See <http://stat.psu.edu/~jiali/course/stat597e/notes2/logit.pdf>

Part III

Distributed Matrix Methods

The Distributed Matrix Data Structure

If I were again beginning my studies, I would follow the advice of Plato and start with mathematics.

—Galileo Galilei

Before continuing, we must spend some time describing a new distributed data structure. In reality, this data structure is the merging of two different kinds of distributed data structures, namely *block distributions* and *cyclic distributions*. Eventually we will get to *block cyclic distributions*, but this structure is complicated enough that it is wise to examine each component separately first.

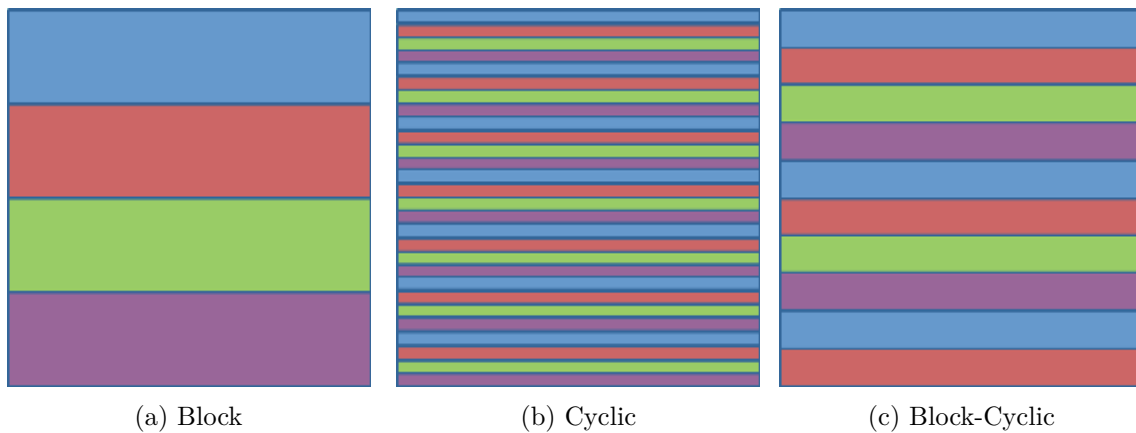


Figure 5.1: Matrix Distribution Schemes

Figure 5.1 shows examples of the three different distribution schemes for 4 processors. The block scheme is simple enough; imagine chopping the matrix into nearly equal blocks and distributing those blocks to different processors. This can be viewed as a special case of the GBD data structure of Section 3.3.

For the cyclic distribution scheme, one can imagine taking each row (or column) of a matrix

and sending the first to one processor, the second to the next, and so on until all processors are exhausted; if the data is not exhausted, then one merely cycles back through the processors, continuing in this fashion until all of the matrix has been distributed.

Finally, the block-cyclic decomposition is the obvious blending of these two schemes, so that each of the former becomes a special case of this new type. Here, we can imagine chopping the matrix up into blocks, but the blocks are not (necessarily) so large that they use up the entire matrix. Once we use up all of the blocks, we use the cyclic data distribution scheme to cycle back through our processors, only using (potentially) blocks of more than one row at a time. From this light, a block-cyclic distribution where the block size is large enough to get all of the data in one cycle is also a block distribution, and a block-cyclic distribution where the blocks are taking just one row at a time is also a cyclic distribution.

The obvious analogue to Figure 5.1 for distributing by column is also possible, but there is a much more important — and complicated — generalization of this scheme. Above, we were thinking of the aggregate of processors as essentially being in a vector, or lying on a one-dimensional line. However, we can extend this to two-dimensional grids of processors as well. Figure 5.2

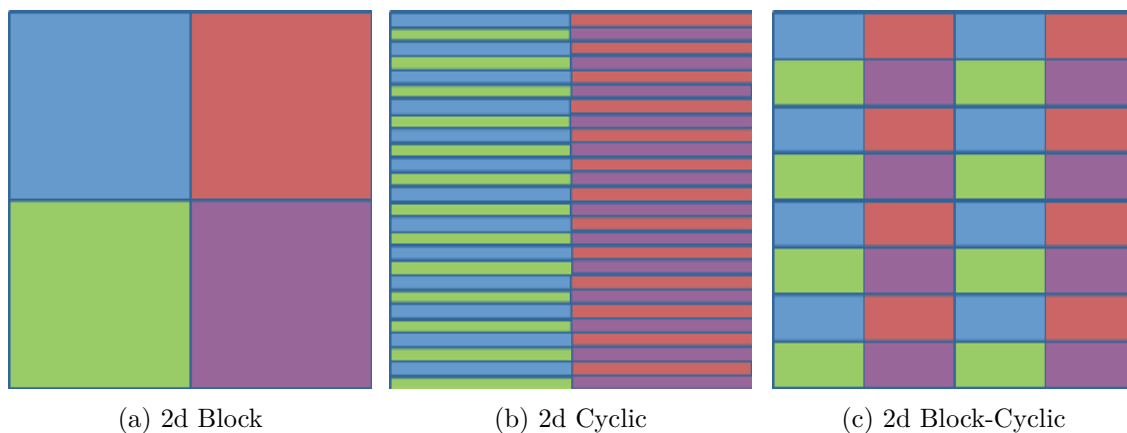


Figure 5.2: Matrix Distribution Schemes Onto a 2-Dimensional Grid

shows how the extension to a 2-dimensional grid of processors, still with just 4 processors, only here, we are assuming that they form a 2×2 grid. This data structure is a generalization of the 1-dimensional block-cyclic distribution, and so it is a generalization of 1-dimensional block and 1-dimensional cyclic distributions as well.

The data structure can get quite complicated, especially when there are many processors involved. Table 5.1 shows the different possible grid shapes for six processors. In general, if we have n processors, then there are $\sigma_0(n)$ total possible grid shapes, where

$$\sigma_m(n) = \sum_{d \mid n} d^m$$

and $d \in \mathbb{N}$ (a positive integer). Thus the grid shapes are given by:

$$\left(d, \frac{n}{d}\right)$$

for each $d \mid n$ with $d \in \mathbb{N}$.

| | | | |
|---|--|---|--|
| $\begin{bmatrix} 0 & 1 & 2 & 3 & 4 & 5 \end{bmatrix}$ | $\begin{bmatrix} 0 & 1 & 2 \\ 3 & 4 & 5 \end{bmatrix}$ | $\begin{bmatrix} 0 & 1 \\ 2 & 3 \\ 4 & 5 \end{bmatrix}$ | $\begin{bmatrix} 0 \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{bmatrix}$ |
| (a) 1×6 | (b) 2×3 | (c) 3×2 | (d) 6×1 |

Table 5.1: Processor Grid Shapes with 6 Processors

This added complication is not for pure masochism; it has some real advantages. For one, this 2-dimensional block-cyclic (henceforth simply referred to as “block-cyclic”) decomposition is the data structure employed by the state of the art dense linear algebra library ScaLAPACK, and if one wishes to use this library, then the use must occur on its terms. However, there are some real performance benefits to this data structure. For many linear algebra operations (which includes many statistical operations, in whole or in part), this data structure offers an interesting balance between communication cost and parallelism. For very large problems, many are surprised to find that communication between processors will often dwarf the computation overhead. This will generally become apparent at the 10,000+ processor count except for the most embarrassingly parallel problems, and the cost of communication gets *much* worse the more cores are added after that. The rate at which this scales badly will depend a great deal on the hardware, but there is no machine in existence at the time of writing for which the above vague warning will not hold true.

Returning to the data structure, notice that since we have control over the processor grid shape and the blocking factor (or blocking dimension — the number of rows/columns in the blocks for the block-decomposition), we can very directly tune the amount of parallelism, and therefore the amount of communication. Make the blocks too small (say 1×1 , or single element blocks) and there will be a great deal of parallelism, in the sense that most processors will stay busy most of the time; but the processors will have to talk to each other to get *anything* done. This makes the communication cost skyrocket. On the other hand, we could make the blocking factor so large in each dimension that it encompasses the entire matrix. That is, the matrix would be stored in its entirety on a single processor. In doing so, we entirely eliminate the communication, but we also eliminate the parallelism.

The fact of the matter is, hard problems require data movement and communication. We should strive to minimize these burdens, but not so myopically that we throw out the parallelism as well. Balancing these parameters then becomes important, and a not entirely trivial optimization problem. The **pbdDMAT** package includes defaults for each that should be “ok” if you have no intuition whatsoever. However, these defaults may not be well-suited to a specific problem, and knowing ahead of time how best to distribute the data is often more art than science.

For the remainder of this chapter, we will be examining these shapes in more depth to get a better feel for the data structure. To do so, let us return to our old friend from Section 3.3:

$$x = \begin{bmatrix} x_{11} & x_{12} & x_{13} & x_{14} & x_{15} & x_{16} & x_{17} & x_{18} & x_{19} \\ x_{21} & x_{22} & x_{23} & x_{24} & x_{25} & x_{26} & x_{27} & x_{28} & x_{29} \\ x_{31} & x_{32} & x_{33} & x_{34} & x_{35} & x_{36} & x_{37} & x_{38} & x_{39} \\ x_{41} & x_{42} & x_{43} & x_{44} & x_{45} & x_{46} & x_{47} & x_{48} & x_{49} \\ x_{51} & x_{52} & x_{53} & x_{54} & x_{55} & x_{56} & x_{57} & x_{58} & x_{59} \\ x_{61} & x_{62} & x_{63} & x_{64} & x_{65} & x_{66} & x_{67} & x_{68} & x_{69} \\ x_{71} & x_{72} & x_{73} & x_{74} & x_{75} & x_{76} & x_{77} & x_{78} & x_{79} \\ x_{81} & x_{82} & x_{83} & x_{84} & x_{85} & x_{86} & x_{87} & x_{88} & x_{89} \\ x_{91} & x_{92} & x_{93} & x_{94} & x_{95} & x_{96} & x_{97} & x_{98} & x_{99} \end{bmatrix}_{9 \times 9}$$

However, we note that the **pbdDMAT** package offers numerous high-level tools for managing these structures, so that the management of distributed details can be as implicit or explicit as the user desires.

5.1 Block Data Distributions

Let us start with the 1-dimensional block data distribution. So here, we will assume that our processor grid looks like:

$$\text{Processors} = \begin{bmatrix} 0 & 1 & 2 & 3 \end{bmatrix}$$

To block-distribute our matrix onto this 1-dimensional grid by rows, then we would have no option but to do the following:

$$x = \begin{bmatrix} x_{11} & x_{12} & x_{13} & x_{14} & x_{15} & x_{16} & x_{17} & x_{18} & x_{19} \\ x_{21} & x_{22} & x_{23} & x_{24} & x_{25} & x_{26} & x_{27} & x_{28} & x_{29} \\ x_{31} & x_{32} & x_{33} & x_{34} & x_{35} & x_{36} & x_{37} & x_{38} & x_{39} \\ \hline x_{41} & x_{42} & x_{43} & x_{44} & x_{45} & x_{46} & x_{47} & x_{48} & x_{49} \\ x_{51} & x_{52} & x_{53} & x_{54} & x_{55} & x_{56} & x_{57} & x_{58} & x_{59} \\ x_{61} & x_{62} & x_{63} & x_{64} & x_{65} & x_{66} & x_{67} & x_{68} & x_{69} \\ \hline x_{71} & x_{72} & x_{73} & x_{74} & x_{75} & x_{76} & x_{77} & x_{78} & x_{79} \\ x_{81} & x_{82} & x_{83} & x_{84} & x_{85} & x_{86} & x_{87} & x_{88} & x_{89} \\ x_{91} & x_{92} & x_{93} & x_{94} & x_{95} & x_{96} & x_{97} & x_{98} & x_{99} \end{bmatrix}_{9 \times 9}$$

Notice that here, processor 3 receives none of the matrix. This is so because if the block size (here, 3×9) were any smaller, then we would not be able to distribute all of the data without cycling. Similarly, if we were to distribute by column then we would have:

$$x = \begin{bmatrix} x_{11} & x_{12} & x_{13} & x_{14} & x_{15} & x_{16} & x_{17} & x_{18} & x_{19} \\ x_{21} & x_{22} & x_{23} & x_{24} & x_{25} & x_{26} & x_{27} & x_{28} & x_{29} \\ x_{31} & x_{32} & x_{33} & x_{34} & x_{35} & x_{36} & x_{37} & x_{38} & x_{39} \\ x_{41} & x_{42} & x_{43} & x_{44} & x_{45} & x_{46} & x_{47} & x_{48} & x_{49} \\ x_{51} & x_{52} & x_{53} & x_{54} & x_{55} & x_{56} & x_{57} & x_{58} & x_{59} \\ x_{61} & x_{62} & x_{63} & x_{64} & x_{65} & x_{66} & x_{67} & x_{68} & x_{69} \\ x_{71} & x_{72} & x_{73} & x_{74} & x_{75} & x_{76} & x_{77} & x_{78} & x_{79} \\ x_{81} & x_{82} & x_{83} & x_{84} & x_{85} & x_{86} & x_{87} & x_{88} & x_{89} \\ x_{91} & x_{92} & x_{93} & x_{94} & x_{95} & x_{96} & x_{97} & x_{98} & x_{99} \end{bmatrix}_{9 \times 9}$$

for exactly the same reason.

If we used a 2-dimensional grid of processors, say a 2×2 grid:

$$\text{Processors} = \begin{bmatrix} 0 & 1 \\ 2 & 3 \end{bmatrix}$$

then our data would be distributed as

$$x = \left[\begin{array}{ccccc|cccc} x_{11} & x_{12} & x_{13} & x_{14} & x_{15} & x_{16} & x_{17} & x_{18} & x_{19} \\ x_{21} & x_{22} & x_{23} & x_{24} & x_{25} & x_{26} & x_{27} & x_{28} & x_{29} \\ x_{31} & x_{32} & x_{33} & x_{34} & x_{35} & x_{36} & x_{37} & x_{38} & x_{39} \\ x_{41} & x_{42} & x_{43} & x_{44} & x_{45} & x_{46} & x_{47} & x_{48} & x_{49} \\ x_{51} & x_{52} & x_{53} & x_{54} & x_{55} & x_{56} & x_{57} & x_{58} & x_{59} \\ \hline x_{61} & x_{62} & x_{63} & x_{64} & x_{65} & x_{66} & x_{67} & x_{68} & x_{69} \\ x_{71} & x_{72} & x_{73} & x_{74} & x_{75} & x_{76} & x_{77} & x_{78} & x_{79} \\ x_{81} & x_{82} & x_{83} & x_{84} & x_{85} & x_{86} & x_{87} & x_{88} & x_{89} \\ x_{91} & x_{92} & x_{93} & x_{94} & x_{95} & x_{96} & x_{97} & x_{98} & x_{99} \end{array} \right]_{9 \times 9}$$

5.2 Cyclic Data Distributions

Proceeding as in the previous section, we would cyclically distribute this matrix by row onto the 1-dimensional processor grid as:

$$x = \left[\begin{array}{ccccccccc} x_{11} & x_{12} & x_{13} & x_{14} & x_{15} & x_{16} & x_{17} & x_{18} & x_{19} \\ \hline x_{21} & x_{22} & x_{23} & x_{24} & x_{25} & x_{26} & x_{27} & x_{28} & x_{29} \\ \hline x_{31} & x_{32} & x_{33} & x_{34} & x_{35} & x_{36} & x_{37} & x_{38} & x_{39} \\ \hline x_{41} & x_{42} & x_{43} & x_{44} & x_{45} & x_{46} & x_{47} & x_{48} & x_{49} \\ \hline x_{51} & x_{52} & x_{53} & x_{54} & x_{55} & x_{56} & x_{57} & x_{58} & x_{59} \\ \hline x_{61} & x_{62} & x_{63} & x_{64} & x_{65} & x_{66} & x_{67} & x_{68} & x_{69} \\ \hline x_{71} & x_{72} & x_{73} & x_{74} & x_{75} & x_{76} & x_{77} & x_{78} & x_{79} \\ \hline x_{81} & x_{82} & x_{83} & x_{84} & x_{85} & x_{86} & x_{87} & x_{88} & x_{89} \\ \hline x_{91} & x_{92} & x_{93} & x_{94} & x_{95} & x_{96} & x_{97} & x_{98} & x_{99} \end{array} \right]_{9 \times 9}$$

and by column:

$$x = \left[\begin{array}{c|c|c|c|c|c|c|c|c} x_{11} & x_{12} & x_{13} & x_{14} & x_{15} & x_{16} & x_{17} & x_{18} & x_{19} \\ \hline x_{21} & x_{22} & x_{23} & x_{24} & x_{25} & x_{26} & x_{27} & x_{28} & x_{29} \\ \hline x_{31} & x_{32} & x_{33} & x_{34} & x_{35} & x_{36} & x_{37} & x_{38} & x_{39} \\ \hline x_{41} & x_{42} & x_{43} & x_{44} & x_{45} & x_{46} & x_{47} & x_{48} & x_{49} \\ \hline x_{51} & x_{52} & x_{53} & x_{54} & x_{55} & x_{56} & x_{57} & x_{58} & x_{59} \\ \hline x_{61} & x_{62} & x_{63} & x_{64} & x_{65} & x_{66} & x_{67} & x_{68} & x_{69} \\ \hline x_{71} & x_{72} & x_{73} & x_{74} & x_{75} & x_{76} & x_{77} & x_{78} & x_{79} \\ \hline x_{81} & x_{82} & x_{83} & x_{84} & x_{85} & x_{86} & x_{87} & x_{88} & x_{89} \\ \hline x_{91} & x_{92} & x_{93} & x_{94} & x_{95} & x_{96} & x_{97} & x_{98} & x_{99} \end{array} \right]_{9 \times 9}$$

Finally, the distribution onto the 2-dimensional grid would look like:

$$x = \left[\begin{array}{ccccc|ccccc} x_{11} & x_{12} & x_{13} & x_{14} & x_{15} & x_{16} & x_{17} & x_{18} & x_{19} \\ x_{21} & x_{22} & x_{23} & x_{24} & x_{25} & x_{26} & x_{27} & x_{28} & x_{29} \\ x_{31} & x_{32} & x_{33} & x_{34} & x_{35} & x_{36} & x_{37} & x_{38} & x_{39} \\ x_{41} & x_{42} & x_{43} & x_{44} & x_{45} & x_{46} & x_{47} & x_{48} & x_{49} \\ x_{51} & x_{52} & x_{53} & x_{54} & x_{55} & x_{56} & x_{57} & x_{58} & x_{59} \\ x_{61} & x_{62} & x_{63} & x_{64} & x_{65} & x_{66} & x_{67} & x_{68} & x_{69} \\ x_{71} & x_{72} & x_{73} & x_{74} & x_{75} & x_{76} & x_{77} & x_{78} & x_{79} \\ x_{81} & x_{82} & x_{83} & x_{84} & x_{85} & x_{86} & x_{87} & x_{88} & x_{89} \\ x_{91} & x_{92} & x_{93} & x_{94} & x_{95} & x_{96} & x_{97} & x_{98} & x_{99} \end{array} \right]_{9 \times 9}$$

5.3 Block-Cyclic Data Distributions

By this time, the reader should feel fairly comfortable with the basic idea and the distribution scheme. So we will jump straight to full generality. To make things more interesting (really, to show the full generality of the distribution), let us now suppose that we have 6 processors in a 2×3 grid:

$$\text{Processors} = \left[\begin{array}{ccc} 0 & 1 & 2 \\ 3 & 4 & 5 \end{array} \right] = \left[\begin{array}{ccc} (0,0) & (0,1) & (0,2) \\ (1,0) & (1,1) & (1,2) \end{array} \right]$$

with the usual MPI processor rank on the left, and the corresponding BLACS processor grid position on the right. This new naming convention is just for convenience of describing a processor by its position in the grid and carries no additional semantic meaning. We will preserve our 2×2 dimensional blocking factor.

Recall that to distribute this data across our 6 processors in the form of a 2×3 process grid in 2×2 blocks, we go in a “round robin” fashion, assigning 2×2 submatrices of the original matrix to the appropriate processor, starting with processor $(0,0)$. Then, if possible, we move on to the next 2×2 block of x and give it to processor $(0,1)$. We continue in this fashion with $(0,2)$ if necessary, and if there is yet more of x in that row still without ownership, we cycle back to processor $(0,0)$ and start over, continuing in this fashion until there is nothing left to distribute in that row.

After all the data in the first two rows of x has been chopped into 2-column blocks and given to the appropriate process in process-column 1, we then move onto the next 2 rows, proceeding in the same way but now using the second process row from our process grid. For the next 2 rows, we cycle back to process row 1. And so on and so forth.

Then distributed across processors, the data will look like:

$$x = \begin{bmatrix} x_{11} & x_{12} & x_{13} & x_{14} & x_{15} & x_{16} & x_{17} & x_{18} & x_{19} \\ x_{21} & x_{22} & x_{23} & x_{24} & x_{25} & x_{26} & x_{27} & x_{28} & x_{29} \\ x_{31} & x_{32} & x_{33} & x_{34} & x_{35} & x_{36} & x_{37} & x_{38} & x_{39} \\ x_{41} & x_{42} & x_{43} & x_{44} & x_{45} & x_{46} & x_{47} & x_{48} & x_{49} \\ x_{51} & x_{52} & x_{53} & x_{54} & x_{55} & x_{56} & x_{57} & x_{58} & x_{59} \\ x_{61} & x_{62} & x_{63} & x_{64} & x_{65} & x_{66} & x_{67} & x_{68} & x_{69} \\ x_{71} & x_{72} & x_{73} & x_{74} & x_{75} & x_{76} & x_{77} & x_{78} & x_{79} \\ x_{81} & x_{82} & x_{83} & x_{84} & x_{85} & x_{86} & x_{87} & x_{88} & x_{89} \\ x_{91} & x_{92} & x_{93} & x_{94} & x_{95} & x_{96} & x_{97} & x_{98} & x_{99} \end{bmatrix}_{9 \times 9}$$

with local storage:

$$\begin{bmatrix} x_{11} & x_{12} & x_{17} & x_{18} \\ x_{21} & x_{22} & x_{27} & x_{28} \\ x_{51} & x_{52} & x_{57} & x_{58} \\ x_{61} & x_{62} & x_{67} & x_{68} \\ x_{91} & x_{92} & x_{97} & x_{98} \end{bmatrix}_{5 \times 4} \begin{bmatrix} x_{13} & x_{14} & x_{19} \\ x_{23} & x_{24} & x_{29} \\ x_{53} & x_{54} & x_{59} \\ x_{63} & x_{64} & x_{69} \\ x_{93} & x_{94} & x_{99} \end{bmatrix}_{5 \times 3} \begin{bmatrix} x_{15} & x_{16} \\ x_{25} & x_{26} \\ x_{55} & x_{56} \\ x_{65} & x_{66} \\ x_{95} & x_{96} \end{bmatrix}_{5 \times 2}$$

$$\begin{bmatrix} x_{31} & x_{32} & x_{37} & x_{38} \\ x_{41} & x_{42} & x_{47} & x_{48} \\ x_{71} & x_{72} & x_{77} & x_{78} \\ x_{81} & x_{82} & x_{87} & x_{88} \end{bmatrix}_{4 \times 4} \begin{bmatrix} x_{33} & x_{34} & x_{39} \\ x_{43} & x_{44} & x_{49} \\ x_{73} & x_{74} & x_{79} \\ x_{83} & x_{84} & x_{89} \end{bmatrix}_{4 \times 3} \begin{bmatrix} x_{35} & x_{36} \\ x_{45} & x_{46} \\ x_{75} & x_{76} \\ x_{85} & x_{86} \end{bmatrix}_{4 \times 2}$$

You *could* use some more natural data distributions than the above, such as the block data structure. However, this may have a substantial impact on performance, depending on the kinds of operations you wish to do. For things that make extensive use of linear algebra — particularly matrix factorizations — you are probably much better off using the above kind of block-cyclic data distribution. Sometimes there is a benefit to using a 1-dimensional grid of processors while still using the full block-cyclic structure. These different processor grid shapes are referred to as *contexts*. They are actually specialized MPI communicators. By default, the recommended (easy) way of managing these contexts with **pbdDMAT** is to call

```
1 library(pbdDMAT, quiet = TRUE)
2 init.grid()
```

The call to `init.grid()` will initialize three such contexts, named 0, 1, and 2. Context 0 is a communicator with processors as close to square as possible, like above. This can be confusing if you ever need to directly manipulate this data structure, but **pbdDMAT** contains *numerous* helper methods to make this process painless, often akin to manipulating an ordinary, non-distributed R data structure. Context 1 puts the processors in a 1-dimensional grid consisting of 1 row. Continuing with our example, the processors form the grid:

$$\text{Processors} = \begin{bmatrix} 0 & 1 & 2 & 3 & 4 & 5 \end{bmatrix} = \begin{bmatrix} (0,0) & (0,1) & (0,2) & (0,3) & (0,4) & (0,5) \end{bmatrix}$$

and if we preserve the 2×2 blocking factor, then the data would be distributed like so:

$$x = \begin{bmatrix} x_{11} & x_{12} & x_{13} & x_{14} & x_{15} & x_{16} & x_{17} & x_{18} & x_{19} \\ x_{21} & x_{22} & x_{23} & x_{24} & x_{25} & x_{26} & x_{27} & x_{28} & x_{29} \\ x_{31} & x_{32} & x_{33} & x_{34} & x_{35} & x_{36} & x_{37} & x_{38} & x_{39} \\ x_{41} & x_{42} & x_{43} & x_{44} & x_{45} & x_{46} & x_{47} & x_{48} & x_{49} \\ x_{51} & x_{52} & x_{53} & x_{54} & x_{55} & x_{56} & x_{57} & x_{58} & x_{59} \\ x_{61} & x_{62} & x_{63} & x_{64} & x_{65} & x_{66} & x_{67} & x_{68} & x_{69} \\ x_{71} & x_{72} & x_{73} & x_{74} & x_{75} & x_{76} & x_{77} & x_{78} & x_{79} \\ x_{81} & x_{82} & x_{83} & x_{84} & x_{85} & x_{86} & x_{87} & x_{88} & x_{89} \\ x_{91} & x_{92} & x_{93} & x_{94} & x_{95} & x_{96} & x_{97} & x_{98} & x_{99} \end{bmatrix}_{9 \times 9}$$

Locally, the data is stored as follows:

$$\begin{bmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \\ x_{31} & x_{32} \\ x_{41} & x_{42} \\ x_{51} & x_{52} \\ x_{61} & x_{62} \\ x_{71} & x_{72} \\ x_{81} & x_{82} \\ x_{91} & x_{92} \end{bmatrix}_{9 \times 2} \begin{bmatrix} x_{13} & x_{14} \\ x_{23} & x_{24} \\ x_{33} & x_{34} \\ x_{43} & x_{44} \\ x_{53} & x_{54} \\ x_{63} & x_{64} \\ x_{73} & x_{74} \\ x_{83} & x_{84} \\ x_{93} & x_{94} \end{bmatrix}_{9 \times 2} \begin{bmatrix} x_{15} & x_{16} \\ x_{25} & x_{26} \\ x_{35} & x_{36} \\ x_{45} & x_{46} \\ x_{55} & x_{56} \\ x_{65} & x_{66} \\ x_{75} & x_{76} \\ x_{85} & x_{86} \\ x_{95} & x_{96} \end{bmatrix}_{9 \times 2} \begin{bmatrix} x_{17} & x_{18} \\ x_{27} & x_{28} \\ x_{37} & x_{38} \\ x_{47} & x_{48} \\ x_{57} & x_{58} \\ x_{67} & x_{68} \\ x_{77} & x_{78} \\ x_{87} & x_{88} \\ x_{97} & x_{98} \end{bmatrix}_{9 \times 2} \begin{bmatrix} x_{19} \\ x_{29} \\ x_{39} \\ x_{49} \\ x_{59} \\ x_{69} \\ x_{79} \\ x_{89} \\ x_{99} \end{bmatrix}_{9 \times 1} \begin{bmatrix} \\ \\ \\ \\ \\ \\ \\ \\ \end{bmatrix}_{0 \times 1}$$

Here, the first dimension of the blocking factor is irrelevant. All processors own either some part of *all* rows, or they own nothing at all. So the above would be the exact same data distribution if we had a blocking factor of 100×2 or 2×2 . However, the decomposition is still block-cyclic; here we use up everything before needing to cycle, based on our choice of blocking factor. If we instead chose a 1×1 blocking, then the data would be distributed like so:

$$x = \begin{bmatrix} x_{11} & x_{12} & x_{13} & x_{14} & x_{15} & x_{16} & x_{17} & x_{18} & x_{19} \\ x_{21} & x_{22} & x_{23} & x_{24} & x_{25} & x_{26} & x_{27} & x_{28} & x_{29} \\ x_{31} & x_{32} & x_{33} & x_{34} & x_{35} & x_{36} & x_{37} & x_{38} & x_{39} \\ x_{41} & x_{42} & x_{43} & x_{44} & x_{45} & x_{46} & x_{47} & x_{48} & x_{49} \\ x_{51} & x_{52} & x_{53} & x_{54} & x_{55} & x_{56} & x_{57} & x_{58} & x_{59} \\ x_{61} & x_{62} & x_{63} & x_{64} & x_{65} & x_{66} & x_{67} & x_{68} & x_{69} \\ x_{71} & x_{72} & x_{73} & x_{74} & x_{75} & x_{76} & x_{77} & x_{78} & x_{79} \\ x_{81} & x_{82} & x_{83} & x_{84} & x_{85} & x_{86} & x_{87} & x_{88} & x_{89} \\ x_{91} & x_{92} & x_{93} & x_{94} & x_{95} & x_{96} & x_{97} & x_{98} & x_{99} \end{bmatrix}_{9 \times 9}$$

Finally, there is context 2. This is deceptively similar to the GBD data structure, but the two are, in general, not comparable. This context puts the processors in a 1-dimensional grid consisting of one column (note the transpose):

$$\text{Processors} = \begin{bmatrix} 0 & 1 & 2 & 3 & 4 & 5 \end{bmatrix}^T = \begin{bmatrix} (0,0) & (1,0) & (2,0) & (3,0) & (4,0) & (5,0) \end{bmatrix}^T$$

So here, the data would be decomposed as:

$$x = \begin{bmatrix} x_{11} & x_{12} & x_{13} & x_{14} & x_{15} & x_{16} & x_{17} & x_{18} & x_{19} \\ x_{21} & x_{22} & x_{23} & x_{24} & x_{25} & x_{26} & x_{27} & x_{28} & x_{29} \\ x_{31} & x_{32} & x_{33} & x_{34} & x_{35} & x_{36} & x_{37} & x_{38} & x_{39} \\ x_{41} & x_{42} & x_{43} & x_{44} & x_{45} & x_{46} & x_{47} & x_{48} & x_{49} \\ x_{51} & x_{52} & x_{53} & x_{54} & x_{55} & x_{56} & x_{57} & x_{58} & x_{59} \\ x_{61} & x_{62} & x_{63} & x_{64} & x_{65} & x_{66} & x_{67} & x_{68} & x_{69} \\ x_{71} & x_{72} & x_{73} & x_{74} & x_{75} & x_{76} & x_{77} & x_{78} & x_{79} \\ x_{81} & x_{82} & x_{83} & x_{84} & x_{85} & x_{86} & x_{87} & x_{88} & x_{89} \\ x_{91} & x_{92} & x_{93} & x_{94} & x_{95} & x_{96} & x_{97} & x_{98} & x_{99} \end{bmatrix}_{9 \times 9}$$

with local storage view:

$$\begin{bmatrix} x_{11} & x_{12} & x_{13} & x_{14} & x_{15} & x_{16} & x_{17} & x_{18} & x_{19} \\ x_{21} & x_{22} & x_{23} & x_{24} & x_{25} & x_{26} & x_{27} & x_{28} & x_{29} \end{bmatrix}_{2 \times 9}$$

$$\begin{bmatrix} x_{31} & x_{32} & x_{33} & x_{34} & x_{35} & x_{36} & x_{37} & x_{38} & x_{39} \\ x_{41} & x_{42} & x_{43} & x_{44} & x_{45} & x_{46} & x_{47} & x_{48} & x_{49} \end{bmatrix}_{2 \times 9}$$

$$\begin{bmatrix} x_{51} & x_{52} & x_{53} & x_{54} & x_{55} & x_{56} & x_{57} & x_{58} & x_{59} \\ x_{61} & x_{62} & x_{63} & x_{64} & x_{65} & x_{66} & x_{67} & x_{68} & x_{69} \end{bmatrix}_{9 \times 2}$$

$$\begin{bmatrix} x_{71} & x_{72} & x_{73} & x_{74} & x_{75} & x_{76} & x_{77} & x_{78} & x_{79} \\ x_{81} & x_{82} & x_{83} & x_{84} & x_{85} & x_{86} & x_{87} & x_{88} & x_{89} \end{bmatrix}_{9 \times 2}$$

$$\begin{bmatrix} x_{91} & x_{92} & x_{93} & x_{94} & x_{95} & x_{96} & x_{97} & x_{98} & x_{99} \end{bmatrix}_{9 \times 1}$$

$$\begin{bmatrix} \end{bmatrix}_{1 \times 0}$$

5.4 Summary

This 2-dimensional block-cyclic data structure — the DMAT data structure — is fairly complicated, but can pay great dividends if some appreciation and understand is given to it. To briefly summarize this data structure:

1. DMAT is *distributed*. No one processor owns all of the matrix.
2. DMAT is *non-overlapping*. Any piece owned by one processor is owned by no other processors.
3. DMAT can be row-contiguous or not, depending on the blocking factor used.
4. Processor 0 = (0,0) will always own at least as much data as every other processor.
5. DMAT is locally column-major and globally, it depends...
6. DMAT is confusing, but very robust and useful for matrix algebra (and thus most non-trivial statistics).

The only items in common between GBD and DMAT are items 1 and 2. A full characterization can be given as follows. Let X be a distributed matrix with n (global) rows and p (global) columns. Suppose we distribute this matrix onto a set of $nprocs$ processors in context 2 using a blocking factor $b = (b_1, b_2)$. Then DMAT is a special case of GBD *if and only if* we have $b_1 > \frac{n}{nprocs}$. Otherwise, there is no relationship between these two structures (and converting between them can be difficult). However, converting between different kinds of block-cyclic layouts is very simple, with numerous high-level methods to assist in this. This process is explained in depth in Section 11.

In the chapters to follow, we offer numerous examples utilizing this data structure. The dedicated reader can find more information about these contexts and utilizing the DMAT data structure, see the **pbdBASE** (Schmidt et al., 2012b) and **pbdDMAT** (Schmidt et al., 2012d) vignettes. Additionally, you can experiment more with different kinds of block-cyclic data distributions on 2-dimensional processor grids using a very useful website at <http://acts.nersc.gov/scalapack/hands-on/datadist.html>.

5.5 Exercises

- 5-1 Experiment with the 2d block-cyclic data layout using this online tool: <http://acts.nersc.gov/scalapack/hands-on/datadist.html> and the **pbdDEMO** function `plot.dmat()`.
- 5-2 Read two papers given at <http://acts.nersc.gov/scalapack/hands-on/datadist.html>. “The Design of Linear Algebra Libraries for High Performance Computers”, by J. Dongarra and D. Walker, and “Parallel Numerical Linear Algebra”, by J. Demmel, M. Heath, and H. van der Vorst.

Constructing Distributed Matrices

*Truth is ever to be found in the simplicity, and
not in the multiplicity and confusion of things.*

—Sir Isaac Newton

The **pbdBASE** and **pbdDMAT** packages offer a distributed matrix class, **ddmatrix**, as well as a collection of high-level methods for performing common matrix operations. For example, if you want to compute the mean of an R matrix **x**, you would call

```
1 mean(x)
```

That’s exactly the same command you would issue if **x** is no longer an ordinary R matrix, but a distributed matrix. These methods range from simple, embarrassingly parallel operations like sums and means, to tightly coupled linear algebra operations like matrix-matrix multiply and singular value decomposition.

Unfortunately, these higher methods come with a different cost: getting the data into the right format, namely the distributed matrix data structure DMAT, discussed at length in the previous chapter. That said, once the hurdle of getting the data into the “right format” is out of the way, these methods offer very simple syntax (designed to mimic R as closely as possible), with the ability to scale computations on very large distributed machines. But to get to the fun stuff, the process of exactly how to decompose data into a block-cyclic distribution must be addressed. We begin dealing with this issue in the simplest way possible.

6.1 Fixed Global Dimension

In these examples, we will examine the case where you know ahead of time what the global number of rows and columns are.

6.1.1 Constructing Simple Distributed Matrices

It is possible to construct fairly simple distributed matrices much in the same way that one can construct simple matrices in R. We can do this using the functions `ddmatrix()` and `as.ddmatrix()`. The former essentially behaves identically to R's own `matrix()` function. This function takes a global input vector/matrix `data=`, as well as the global number of rows `nrow=` and the global number of columns `ncol=`. Additionally, the user may specify the blocking factor `bldim=` and the BLACS context `CTXT`, and the return is a distributed matrix. For instance, we can specify

```
ddmatrix()
1 dx <- ddmatrix(data=0, nrow=10, ncol=10)
```

to get a distributed matrix with *global* dimension 10×10 consisting of zeros. We can also do cute things like

```
ddmatrix()
1 dx <- ddmatrix(data=1:3, nrow=5, ncol=5)
```

which will create the distributed analogue of

| | [,1] | [,2] | [,3] | [,4] | [,5] |
|------|------|------|------|------|------|
| [1,] | 1 | 3 | 2 | 1 | 3 |
| [2,] | 2 | 1 | 3 | 2 | 1 |
| [3,] | 3 | 2 | 1 | 3 | 2 |
| [4,] | 1 | 3 | 2 | 1 | 3 |
| [5,] | 2 | 1 | 3 | 2 | 1 |

How exactly that “distributed analogue” will look (locally) depends on the processor grid shape (whence too, the number of processors) as well as the blocking factor. This operation performs no communication.

While this can be useful, it is far from the only way to construct distributed matrices. One can also convert a global (non-distributed) matrix into a distributed matrix. There are some caveats; this matrix must either be owned in total by all processors (which is very useful in testing, but should not be used at scale), or the matrix is owned in total by one processor, with all others owning NULL for that object.

For example, we can create identical return to the above via

```
as.ddmatrix()
1 x <- matrix(data=1:3, nrow=5, ncol=5)
2 dx <- as.ddmatrix(x)
```

or

```
as.ddmatrix()
1 if (comm.rank()==0) {
```

```

2  x <- matrix(data=1:3, nrow=5, ncol=5)
3  } else {
4    x <- NULL
5  }
6
7  dx <- as.ddmatrix(x)

```

Each of these operations performs communication.

Other, more general combinations are possible through other means, but they are much more cumbersome.

6.1.2 Diagonal Distributed Matrices

*Example: construct **diagonal** distributed matrices of specified global dimension.*

The demo command is

Shell Command

```

### At the shell prompt, run the demo with 4 processors by
### (Use Rscript.exe for windows system)
mpiexec -np 4 Rscript -e
    "demo(randmat_diag_global,'pbdDEMO',ask=F,echo=F)"

```

In R, the `diag()` function serves two purposes; namely, it is both a reduction operation and a reverse-reduction operation, depending on the input. More specifically, if given a matrix, it produces a vector containing the diagonal entries of that matrix; but if given a vector, it constructs a diagonal matrix whose diagonal is that vector. And so for example, the zero and identity matrices of any dimension can quickly be constructed via:

Diagonal Matrices in R

```

1  diag(x=0, nrow=10, ncol=10) # zero matrix
2  diag(x=1, nrow=10, ncol=10) # identity matrix

```

Both of the above functionalities of `diag()` are available for distributed matrices; however we will only focus on the latter.

When you wish to construct a diagonal distributed matrix, you can easily do so by using the additional `type=` argument to our `diag()` method. By default, `type="matrix"`, though the user may specify `type="ddmatrix"`. If so, then as one might expect, the optional `bldim=` and `ICTXT=` arguments are available. So with just a little bit of tweaking, the above example becomes:

Diagonal Matrices in pbdR

```

1  diag(x=0, nrow=10, ncol=10, type="ddmatrix") # zero
    (distributed) matrix
2  diag(x=1, nrow=10, ncol=10, type="ddmatrix") # identity
    (distributed) matrix

```

In fact, the `type=` argument employs partial matching, so if we really want to be lazy, then we could simply do the following:

Diagonal Matrices in pbdR

```
1 diag(x=0, nrow=10, ncol=10, type="d") # zero (distributed) matrix
2 diag(x=1, nrow=10, ncol=10, type="d") # identity (distributed)
  matrix
```

Beyond the above brief explanation, the demo for this functionality is mostly self-contained, although we do employ the `redistribute()` function to fully show off local data storage. This function is explained in detail in Chapter 11.

6.1.3 Random Matrices

Example: randomly generate distributed matrices with random normal data of specified global dimension.

The demo command is

Shell Command

```
### At the shell prompt, run the demo with 4 processors by
### (Use Rscript.exe for windows system)
mpirun -np 4 Rscript -e "demo(randmat_global, 'pbdDEMO', ask=F, echo=F)"
```

This demo shows 3 separate ways that one can generate a random normal matrix with specified global dimension. The first two generate the matrix in full on at least one processor and distribute(s) the data, while the last method generates locally only what is needed. As such, the first two can be considered demonstrations with what to do when you have data read in on one processor and need to distribute it out to the remaining processors, but for the purposes of building a randomly generated distributed matrix, they are not particularly efficient strategies.

As described in the previous section, if we have a matrix `x` stored on processor 0 and `NULL` on the others, then we can distribute it out as an object of class `ddmatrix` via the command `as.ddmatrix()`. For example

```
1 if (comm.rank()==0){
2   x <- matrix(rnorm(100), nrow=10, ncol=10)
3 } else {
4   x <- NULL
5 }
6
7 dx <- as.ddmatrix(x)
```

will distribute the required data to the remaining processors. We note for clarity that this is not equivalent to sending the full matrix to all processors and then throwing away all but what is needed. Only the required data is communicated to the processors.

That said, having all of the data on all processors can be convenient while testing, if only for being more minimalistic in the amount of code/thinking required. To do this, one need only do the following:

```
1 x <- matrix(rnorm(100), nrow=10, ncol=10)
2
3 dx <- as.ddmatrix(x)
```

Here, each processor generates the full, global matrix, then throws away what is not needed. Again, this is not efficient, but the code is concise, so it is extremely useful in testing. Now, this assumes you are using the same seed on each processor. This can be managed using the **pbdMPI** function `comm.set.seed()`, as in the demo script. For more information, see that package's documentation.

Finally, you can generate locally only what you need. The demo script does this via the **pbdDMAT** package's `ddmatrix()` function. This is “new” behavior for this syntax (if you view `ddmatrix()` as an extension of `matrix()`). Ordinarily you would merely execute something like

Creating a random normal matrix in serial R

```
1 x <- rnorm(n*p)
2 x <- matrix(x, nrow=n, ncol=p) # this creates a copy
3
4 y <- rnorm(n*p)
5 dim(y) <- c(n, p) # this does not
```

However, things are slightly more complicated with `ddmatrix` objects, and the user may not easily know ahead of time what the size of the local piece is just from knowing the global dimension. Because this requires a much stronger working knowledge of the underlying data structure than most will be comfortable with, we provide this simple functionality as an extension. However, we note that the disciplined reader is more than capable of figuring out how it functions by examining the source code and checking with the reference manual. the size of the local storage. This is all very well documented in the **pbdBASE** documentation, but since no one even pretends to read that stuff, **NUMROC** is a ScaLAPACK tool, which means “NUMber of Rows Or Columns.” The function `base.numroc()` is an implementation in R which calculates the number of rows *and* columns at the same time (so it is a bit of a misnomer, but preserved for historical reasons). dimension `dim`, a blocking factor `bldim`, and a BLACS context number `ICTXT`. The extra argument `fixme` determines whether or not the lowest value returned should be 1. If `fixme==FALSE` and any of the returned local dimensions are less than 1, then that processor does not actually own any of the global matrix — it has no local storage. But something must be stored, and so we default this to `matrix(0)`, the 1×1 matrix with single entry 0.

6.2 Fixed Local Dimension

Example: randomly generate distributed matrices with random normal data of specified local dimension.

The demo command is

Shell Command

```
### At the shell prompt, run the demo with 4 processors by  
### (Use Rscript.exe for windows system)  
mpiexec -np 4 Rscript -e "demo(randmat_local,'pbdDEMO',ask=F,echo=F)"
```

This is similar to the above, but with a critical difference. Instead of specifying a fixed *global* dimension and then go determine what the local storage space is, instead we specify a fixed *local* dimension and then go figure out what the global dimension should be. This can be useful for testing weak scaling of an algorithm, where different numbers of cores are used with the same local problem size.

To this end, the demo script utilizes the `ddmatrix.local()` function, which has the user specify a local dimension size that all the processors should use, as well as a blocking factor and BLACS context. Now here things get somewhat tricky, because in order for this matrix to exist at all, each margin of the blocking factor must divide (as an integer) the corresponding margin of the global dimension. To better understand why this is so, the reader is suggested to read the **pbdDMAT** vignette. But if you already understand or are merely willing to take it on faith, then you surely grant that this is a problem.

So here, we assume that the local dimension is chosen appropriately by the user, but it is possible that a bad blocking factor is supplied by the user. Rather than halt with a stop error, we attempt to find the next best blocking factor possible. To do this, we must find the smallest integer above the specified blocking factor that will divide the number of local rows or columns.

6.3 Exercises

- 6-1 Random number generation (RNG) is used in this Section such as `rnorm()`. In **pbdR**, we use an R package **rlecuyer** (Sevcikova and Rossini, 2012) to set different streams of seed in parallel. Try to find and use other RNG methods or implementations in R.

Basic Examples

I must meditate further on this

—Joseph Louis Lagrange

There is a deep part of the author that does not want to begin with these examples. There is a real danger for the cursory observer to see these and hastily conclude that our work, or R as a whole, is merely a “Matlab Clone.” Nothing could be further from reality.

Matlab is an amazing product. It costs quite a lot of money; it had better damn well be. However, for statistics, machine learning, data mining — data science — we believe that R is “better.” Is R faster? Emphatically, no. But we argue that R wins in other ways.

It is true that everything R can do, so too can Matlab; of course, the converse is also true — that everything Matlab can do, R can do as well. Each is a turing complete language. But being turing complete is not sufficient; \LaTeX is turing complete, and yet we do not perform scientific computation in it (although of course it is unparalleled in typesetting). But we could.

The fact that we do not is an extension of the fact that math journals do not publish articles written in C or Fortran. Those programming languages are the wrong mediums of abstraction for expressing highly complicated ideas to domain experts. Only a madman would attempt to express deep mathematical abstraction in these languages for publication (implementation being an entirely separate issue). Likewise, we do not perform our statistical analyses in \LaTeX (don’t be a pedant; we are not talking about sweave and you know it). People overwhelmingly choose R for the analysis of data because it is the closest brain \rightarrow computer translation available for such problems.

Of course, this goes both ways. If your life is matrix algebra, then R is a much worse fit for you than is Matlab. Much of statistics is applied matrix algebra, but not all matrix algebra is statistics.

So we reluctantly press on with several basic examples utilizing distributed matrices. For meatier examples, see Chapter 8.

7.1 Reductions and Transformations

7.1.1 Reductions

In Section 6.1.2, we discussed the way that the `diag()` method may be utilized as a reduction operator. We have numerous other reductions available, such as `sum()`, `prod()`, `min()`, and `max()`. These operate exactly as their serial counterparts:

Reductions

```
1 library(pbdDMAT, quiet = TRUE)
2 init.grid()
3
4 dx <- ddmatrix(data=0:1, nrow=10, ncol=10)
5
6 sm <- sum(dx)
7 comm.print(sm)
8
9 pd <- prod(dx)
10 comm.print(pd)
11
12 mn <- min(dx)
13 comm.print(mn)
14
15 mx <- max(dx)
16 comm.print(mx)
17
18 finalize()
```

We also offer some “super reductions”. It is possible to change a distributed matrix into a non-distributed matrix or vector using the methods `as.matrix()` or `as.vector()`. For example:

Super Reductions

```
1 library(pbdDMAT, quiet = TRUE)
2 init.grid()
3
4 dx <- ddmatrix(data=0:1, nrow=10, ncol=10)
5 print(dx)
6
7 x <- as.matrix(dx)
8 comm.print(x)
9
10 finalize()
```

These can be very useful in testing, but should be used sparingly at scale.

7.1.2 Transformations

We also offer numerous in-place transformations, such as the various `log()` functions, `abs()`, `sqrt()`, `ceiling()`, `floor()`, and `round()`. For example:

Transformations

```
1 library(pbdDMAT, quiet = TRUE)
2 init.grid()
3
4 comm.set.seed(1234, diff = TRUE)
5
6 dx <- ddmatrix(data=-3:3, nrow=10, ncol=10)
7
8 dx <- ceiling(sqrt(abs(dx)))
9
10 x <- as.matrix(dx)
11 comm.print(x)
12
13 finalize()
```

7.2 Matrix Arithmetic

We also offer a complete set of methods for distributed matrix arithmetic. With identical syntax to R, we can do some reasonably complicated things, such as:

Transformations

```
1 library(pbdDMAT, quiet = TRUE)
2 init.grid()
3
4 dx <- ddmatrix(data=-3:3, nrow=10, ncol=10)
5 vec <- 1:2
6
7 dy <- (dx - vec) %*% dx
8
9 y <- as.matrix(dy)
10 comm.print(y)
11
12 finalize()
```

For a full list of methods, see the **pbdDMAT** documentation.

One item worth noting is that, as with regular R, if the user wishes to compute $X^T X$ or XX^T , then it is usually much faster to use the methods `crossprod()` and `tcrossprod()`, respectively. However, for this operation, things are somewhat more complicated in the distributed sphere

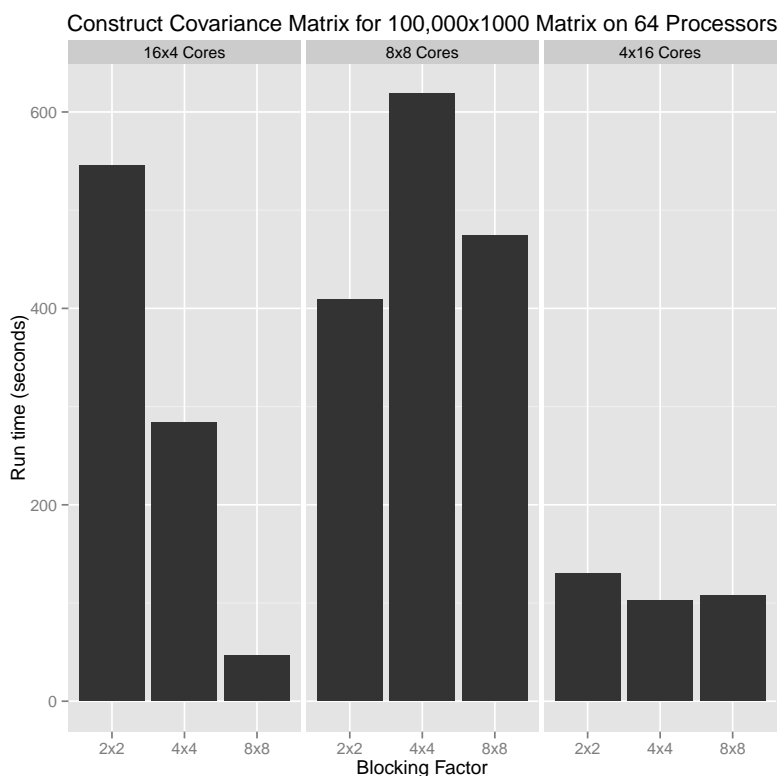


Figure 7.1: Covariance Benchmark Showing Effect of Parameter Calibration

than in serial. Figure 7.1 shows the results of a benchmark of the `cov()` method for computing variance-covariance matrices (which is just a small amount of extra work on top of `crossprod()`). Here, each run consists of 25 replicates of calling `cov()` (which calls `crossprod()`) and then reporting the average run time. The changes in parameters are subtle, but the effects are enormous. Sometimes it may be (much) more beneficial to use `t(x) %*% x`. Others it may not. Proper calibration of these parameters to achieve optimal performance for a given task is still somewhat of an open question to the HPC community.

7.3 Matrix Factorizations

In addition to all of the above, we also provide several of the more important matrix factorizations for distributed matrices. Namely, the singular value decomposition `svd()/La.svd()`, QR factorization `qr()`, Cholesky factorization `chol()`, and LU factorization `lu()`. So for example:

Matrix Factorizations

```

1 library(pbdDEMO, quiet = TRUE)
2 init.grid()
3
4 comm.set.seed(1234, diff = TRUE)
5
```

```

6 dx <- ddmatrix("rnorm", nrow=10, ncol=10, bldim=2)
7
8 out <- chol(crossprod(dx))
9 print(out)
10
11 finalize()

```

7.4 Exercises

7-1 Sub-setting, selection and filtering are basic matrix operation featured in R. The next may look silly, but it is useful for data processing. Suppose \mathbf{X} is in `ddmatrix` with dimension 97×42 , say `dx <- ddmatrix(rnorm(97 * 42), nrow=37)`, do the following:

- `dY <- dx[c(1, 41, 5, 4, 3),]`
`dY <- dx[, c(10:3, 5, 5)]`
`dY <- dx[3:51, 5:3]`
- `dY <- dx[dx[, 31] > 10,]`
`dY <- dx[dx[, 41] > dx[, 40],]`
`dY <- dx[, dx[41,] > dx[40,]]`
`dY <- dx[dx[, 41] > dx[, 40], c(1, 3, 5)]`
- `dx[c(1, 41, 5, 4, 3),] <- 10`
`dx[, c(10:3, 5, 5)] <- 9`
`dx[3:51, 5:3] <- 8`
- `dx[dx[, 31] > 0,] <- 7`
`dx[dx[, 41] > dx[, 40],] <- 6`
`dx[, dx[41,] > dx[40,]] <- 5`
`dx[dx[, 41] > dx[, 40], c(1, 3, 5)] <- 4`
- `dx[c(1, 40, 5, 4, 3),] <- dx[c(1, 41, 5, 4, 3) + 1,]`
`dx[, c(10:3, 5, 5)] <- dx[, c(10:3, 5, 5) + 1]`
`dx[c(10:3, 5, 5),] <- dx[c(10:3, 5, 5) + 1,]`
`dx[3:51, 5:3] <- dx[(3:51) + 1, (5:3) + 1]`
- `dx[dx[, 31] > 0,] <- dx[dx[, 31] > 0, c(42, 1:41)]`
`dx[dx[, 41] > dx[, 40],] <- dx[dx[, 41] > dx[, 40], c(41:42, 1:40)]`
`dx[, dx[41,] > dx[40,]] <- dx[c(96:97, 1:95), dx[, 41] > dx[, 40]]`
`dx[dx[, 41] > dx[, 40], c(1, 3, 5)] <- dx[dx[, 41] > dx[, 40], c(1, 3, 5) + 1]`

If any of above does not work, please report the bugs.

7-2 Suppose `dx` is as Exercise 7-1, do the following:

- `dY <- dx[-c(1, 41, 5, 4, 3),]`
`dY <- dx[, -c(10:3, 5, 5)]`

- ```
dY <- dX[-(3:51), -(5:3)]
```
- `dY <- dX[dX[, 41] > dx[, 40], -c(1, 3, 5)]`
  - `dX[-c(1, 41, 5, 4, 3),] <- 10`  
`dX[, -c(10:3, 5, 5)] <- 9`  
`dX[-(3:51), -(5:3)] <- 8`
  - `dX[dX[, 41] > dx[, 40], -c(1, 3, 5)] <- 4`
  - `dX[-c(1, 40, 5, 4, 3),] <- dX[-(c(1, 41, 5, 4, 3) + 1),]`  
`dX[, -c(10:3, 5, 5)] <- dX[, -(c(10:3, 5, 5) + 1)]`  
`dX[-c(10:3, 5, 5),] <- dX[-(c(10:3, 5, 5) + 1),]`  
`dX[-(3:51), -(5:3)] <- dX[-((3:51) + 1), -((5:3) + 1)]`
  - `dX[dX[, 41] > dx[, 40], -c(1, 3, 5)] <- dX[dX[, 41] > dx[, 40], -(c(1, 3, 5) + 1)]`

7-3 Verify the validity of Exercises 7-1 and 7-2 using ordinary R operations (cast the matrix as global first using `X <- as.matrix(dX)`).

7-4 Implement GBD row-major matrix format in 2 processors for Exercises 7-1 and 7-2.

## Advanced Statistics Examples

*I see it, but I don't believe it.*

—Georg Cantor

The **pbdDMAT** package contains many useful methods for doing computations with distributed matrices. For comprehensive (but shallow) demonstrations of the distributed matrix methods available, see the **pbdDMAT** package's vignette and demos.

Here we showcase a few more advanced things that can be done by chaining together R and **pbdR** code seamlessly.

### 8.1 Sample Mean and Variance Revisited

*Example: Get summary statistics from a distributed matrix.*

The demo command is

Shell Command

```
At the shell prompt, run the demo with 4 processors by
(Use Rscript.exe for windows system)
mpiexec -np 4 Rscript -e "demo(sample_stat_dmat, 'pbdDEMO', ask=F, echo=F)"
```

Returning to the sample statistics problem from Section 4.2, we can solve these same problems — and then some — using distributed matrices. For the remainder, suppose we have a distributed matrix **dx**.

Computing a mean is simple enough. We need only call

Summary Statistics

```
1 mean(dx)
```

We also have access to the other summary statistics methods for matrices, however, such as `rowSums()`, `rowMeans()`, etc. Furthermore, we can calculate variances for distributed matrices. Constructing the variance-covariance matrix is as simple as calling

## Summary Statistics

```
1 cov(dx)
```

Or we could generate standard deviations column-wise, using the method R suggests for ordinary matrices using `apply()`

## Summary Statistics

```
1 apply(dx, MARGIN=2, FUN=sd)
```

or we could simply call

## Summary Statistics

```
1 sd(dx)
```

In R, calling `sd()` on a matrix issues a warning, telling the user to instead use `apply()`. Either of these approaches works with a distributed matrix (with the code as above), but for us, using `sd()` is preferred. This is because, as outlined in Section 11.2, our `apply()` method carries an implicit data redistribution with it, while the `sd()` method is fast, ad-hoc code which requires no redistribution of the data.

## 8.2 Verification of Distributed System Solving

*Example: Solve a system of equations and verify that the solution is correct.*

The demo command is

## Shell Command

```
At the shell prompt, run the demo with 4 processors by
(Use Rscript.exe for windows system)
mpexec -np 4 Rscript -e "demo(verify,'pbdDEMO',ask=F,echo=F)"
```

The **pbdDEMO** contains a set of verification routines, designed to test for validity of the numerical methods at any scale. Herein we will discuss the verification method for solving systems of linear equations, `verify.solve()`.

The process is simple. The goal is to solve the equation (in matrix notation)

$$Ax = b$$

for  $n \times n$  matrix  $A$  and  $n \times 1$  matrix  $b$ . However, here we start with  $A$  and  $x$  and use these to produce  $b$ . We then forget we ever knew what  $x$  was and solve the system. Finally, we remember what  $x$  really should be and compare that with our numerical solution.

More specifically, we take the matrix  $A$  to be random normal generated data and the true solution  $x$  to be a constant vector. We then calculate

$$b := Ax$$

and finally the system is solve for a now (pretend) unknown  $x$ , so that we can compare the numerically determined  $x$  to the true constant  $x$ . All processes are timed, and both success/failure and timing results are printed for the user at the completion of the routine. This effectively amounts to calling:

#### Verifying Distributed System Solving

```

1 # generating data
2 timer({
3 x <- ddmatrix("rnorm", nrow=nrows, ncol=ncols)
4 truesol <- ddmatrix(0.0, nrow=nrows, ncol=1)
5 })
6
7 timer({
8 rhs <- x %*% truesol
9 })
10
11 # solving
12 timer({
13 sol <- solve(x, rhs)
14 })
15
16 # verifying
17 timer({
18 iseq <- all.equal(sol, truesol)
19 iseq <- as.logical(allreduce(iseq, op='min'))
20 })

```

with some added window dressing.

### 8.3 Compression with Principal Components Analysis

*Example: Take PCA and retain only a subset of the rotated data.*

The demo command is

#### Shell Command

```

At the shell prompt, run the demo with 4 processors by
(Use Rscript.exe for windows system)
mpiexec -np 4 Rscript -e "demo(pca,'pbdDEMO',ask=F,echo=F)"

```

Suppose we wish to perform a principal components analysis and retain only some subset of the columns of the rotated data. One of the ways this is often done is by using the singular



values — the standard deviations of the components — as a measure of variation retained by a component. However, the first step is to get the principal components data. Luckily this is trivial. If our data is stored in the distributed matrix object `dx`, then all we need to is issue the command:

```
1 pca <- prcomp(x=dx, retx=TRUE, scale=TRUE)
```

Now that we have our PCA object (which has the same structure as that which comes from calling `prcomp()` on an ordinary R matrix), we need only decide how best to throw away what we do not want. We might want to retain at least as many columns as would be needed to retain 90% of the variation of the original data:

```
1 prop_var <- cumsum(pca$sdev)/sum(pca$sdev)
2 i <- min(which(prop_var > 0.9))
3
4 new_dx <- pca$x[, 1:i]
```

Or we might want one fewer column than the number that would give us 90%:

```
1 prop_var <- cumsum(pca$sdev)/sum(pca$sdev)
2 i <- max(min(which(prop_var > 0.9)) - 1, 1)
3
4 new_dx <- pca$x[, 1:i]
```

## 8.4 Predictions with Linear Regression

*Example: Fit a linear regression model and use it to make a prediction on new data.*

The demo command is

### Shell Command

```
At the shell prompt, run the demo with 4 processors by
(Use Rscript.exe for windows system)
mpiexec -np 4 Rscript -e "demo(ols_dmat,'pbdDEMO',ask=F,echo=F)"
```

Suppose we have some predictor variables stored in the distributed  $n \times p$  matrix `dx` and a response variable stored in the  $n \times 1$  distributed matrix `dy`, and we wish to use the ordinary least squares model from (4.7) to make a prediction about some new data, say `dy.new`. Then this really amounts to just a few simple commands, namely:

```
1 mdl <- lm.fit(dx, dy)
2
3 pred <- dx.new %*% mdl$coefficients
4
5 comm.print(submatrix(pred), quiet=T)
```

## 8.5 Exercises

- 8-1 Based on Section 8.2, extend the code to find  $\mathbf{X}$  which solves  $\mathbf{AX} = \mathbf{B}$  where  $\mathbf{A}$ ,  $\mathbf{X}$  and  $\mathbf{B}$  are matrices with appropriated dimensions and  $\mathbf{A}$  and  $\mathbf{B}$  are known.
- 8-2 The `prcomp()` method introduced in Section 8.3 also returns rotations for all components. Computationally verify with several examples that these rotations are orthogonal, i.e., that their crossproduct is the identity matrix.
- 8-3 Based on Section 8.4, find a point-wise 95% confidence interval for the observed data  $\hat{\mathbf{y}}|\mathbf{X}$  and a 95% predictive interval for the prediction for a new data  $\hat{\mathbf{y}}_{new}|\mathbf{x}_{new}$ .

## Part IV

# Reading and Managing Data

## Reading CSV and SQL Files

*“Data! Data! Data!” he cried impatiently. “I can’t make bricks without clay.”*

—Sherlock Holmes

As we mentioned at the beginning of the discussion on distributed matrix methods, most of the hard work in using these tools is getting the data into the right format. Once this hurdle has been overcome, the syntax will magically begin to look like native R syntax. Some insights into this difficulty can be seen in the previous section, but now we tackle the problem head on: how do you get real data into the distributed matrix format?

### 9.1 CSV Files

*Example: Read data from a csv directly into a distributed matrix.*

The demo command is

#### Shell Command

```
At the shell prompt, run the demo with 4 processors by
(Use Rscript.exe for windows system)
mpiexec -np 4 Rscript -e "demo(read_csv,'pbdDEMO',ask=F,echo=F)"
```

It is simple enough to read in a csv file serially and then distribute the data out to the other processors. This process is essentially identical to one of the random generation methods in Section 6.1.3. For the sake of completeness, we present a simple example here:

```
1 if (comm.rank()==0){ # only read on process 0
2 x <- read.csv("myfile.csv")
3 } else {
4 x <- NULL
```

```
5 }
6
7 dx <- as.ddmatrix(x)
```

However, this is inefficient, especially if the user has access to a parallel file system. In this case, several processes should be used to read parts of the file, and then distribute that data out to the larger process grid. Although really, the user should not be using `csv` to store large amounts of data because it always requires a sort of inherent “serialness”. Regardless, a demonstration of how this is done is useful. We can do so via the **pbdDEMO** package’s function `read.csv.ddmatrix` on an included dataset:

#### Reading a CSV with Multiple Readers

```
1 dx <- read.csv.ddmatrix("../extra/data/x.csv",
2 sep=";", nrows=10, ncols=10,
3 header=TRUE, bldim=4,
4 num.rdrs=2, ICTXT=0)
5
6 print(dx)
```

The code powering the function itself is quite complicated, going well beyond the scope of this document. To understand it, the reader should see the advanced sections of the **pbdDMAT** vignette.

## 9.2 Exercises

- 9-1 In Section 9.1, we have seen an CSV reading example, however, this is not an efficient way for large CSV files by calling `read.csv()`. The R functions `con <- file(...)` can open a connection to the CSV files and `readLines(con, n = 100000)` can access a chunk of data (100,000 lines) from disk more efficiently. Implement a simple function as `read.csv()` and compare performance.
- 9-2 As Exercise 9-1, implement a simple function by utilizing `writeLines()` for writing large CSV file and compare performance to the `write.csv()`.
- 9-3 **pbdMPI** since version 0.2-2 has new functions for simple data input and output (I/O) that functions `comm.read*()` and `comm.write*()` can do either serial or parallel I/O to and from text or csv files. Modify the example of Section 9.1 and compare performance from the above Exercise with those functions in **pbdMPI**.
- 9-4 Other R packages can deal with fast reading for CSV format or so in serial. Try **ff** (Adler et al., 2013) and **bigmemory** (Kane and Emerson, 2010).

*I don't believe in natural science.*

—Kurt Gödel

## 10.1 Introduction

Network Common Data Form version 4 (NetCDF4) is a self-describing, machine-independent data format primarily used for very large scale array-oriented scientific data. The NetCDF4 library is available from the Unidata Program at <http://www.unidata.ucar.edu/software/netcdf>. NetCDF4 is built on top of HDF5 data model for extremely large and complex data collections. More specifically, NetCDF4 is a subset of HDF5 but with enhanced usability features. The HDF5 library is available from the HDF Group <http://www.hdfgroup.org/HDF5/>.

Both libraries provide high-performance functionality to create, access, read, write, and modify NetCDF4 files. The R package **ncdf4** (Pierce, 2012) provides an R-level interface for NetCDF4 libraries. A short summary of its major functions is given in the Table 10.1

Both NetCDF4 and HDF5 provide the capability for parallel I/O, allowing multiple processors to collectively access the same file. To enable this mechanism, HDF5 and NetCDF4 are required to be compiled and linked against an MPI library. In addition to offering access to collective I/O supported by parallel HDF5 and NetCDF4 libraries, the R package **pbdNCDF4** (Patel et al., 2013a) is a parallel extension of **ncdf4** and provides functions for collectively accessing the same NetCDF4 file by multiple processors at the same time.

Users are encouraged to read the vignette (Patel et al., 2013b) of **pbdNCDF4** which includes information for compiling HDF5 and NetCDF4 in parallel, and demonstration of parallel-enabled functions. Table 10.1 also lists the the major functions of **pbdNCDF4**.

The **pbdDEMO** has an example dataset **TREFHT** from a Community Atmosphere Model (CAM) version 5 simulation output. CAM is a series of global atmosphere models originally developed at the National Center for Atmospheric Research (NCAR) and currently guided by Atmosphere Model Working Group (AMWG) of the Community Earth System Model (CESM) project.

Table 10.1: Functions from **pbdNCDF4** and **ncdf4** for accessing NetCDF4 files

| Package                              | Function                       | Purpose                           |
|--------------------------------------|--------------------------------|-----------------------------------|
| <b>pbdNCDF4</b>                      | <code>nc_create_par</code>     | Create a NetCDF4 file in parallel |
|                                      | <code>nc_open_par</code>       | Open a NetCDF4 file in parallel   |
|                                      | <code>nc_var_par_access</code> | Specify parallel variable         |
| <b>ncdf4</b>                         | <code>nc_create</code>         | Create a NetCDF4 file             |
|                                      | <code>nc_open</code>           | Open a NetCDF4 file               |
| <b>pbdNCDF4</b><br>&<br><b>ncdf4</b> | <code>ncdim_def</code>         | Define data dimension             |
|                                      | <code>ncvar_def</code>         | Define a variable                 |
|                                      | <code>ncvar_put</code>         | Write data to a NetCDF4 file      |
|                                      | <code>ncvar_get</code>         | Read data from a NetCDF4 file     |
|                                      | <code>nc_close</code>          | Close a NetCDF4 file              |

CAM version 5 (CAM5) is the latest standalone model modified substantially with a range of enhancements and improvement in the representation of physical processes since version 4 (Eaton, 2011; Vertenstein et al., 2011).

The data TREFHT as shown in the Figure 10.1 is taken from monthly averaged temperature at reference height of January 2004. This dataset is about three megabytes and is a tiny part of ultra-large simulations conducted by Prabhat and Michael Wehner of Lawrence Berkeley National Laboratory. The simulations run from 1987 to 2005 over 1152 longitudes (lon), 768 latitudes (lat), and 30 altitudes (lev). The total amount of simulation outputs is over 200 Terabytes, which are summarized and averaged including monthly-averaged, daily-averaged, and three-hours-averaged data. More datasets are available on ESGF (<http://www.earthsystemgrid.org/>) through the C20C project (on the NERSC portal).

A user with **pbdDEMO** installed can load the TREFHT dataset in the usual way, namely `data(TREFHT)` after loading the **pbdDEMO** package. Here, `TREFHT` is a list consisting of several elements. First, `TREFHT$def` contains all definitions regarding to this variable in class `ncvar4` including locations, dimensions, units, variable size, data storage, missing values, etc.

Next, `TREFHT$def$size` gives the data dimensions which are  $(lon, lat, time) = (1152, 768, 1)$ . Since this data is monthly averaged of Jan. 2004, it is stored as an one-time step output which is an averaged slice among 20 years.

Finally, `TREFHT$data` contains the values of each location and is a matrix with dimension  $1152 \times 768$ . Note that the column (lon) is in x-axis direction and the row (lat) is in y-axis direction.

*Example: Temperature at reference height (TREFHT).*

In an R session (interactive mode), run the demo by executing

R Code

```
demo(trefht, 'pbdDEMO', ask = F, echo = F)
```

This will show a plot as the Figure 10.1 providing a visualization about this variable and how temperatures are vary across locations, particularly decreasing in latitudes. Moreover, the South

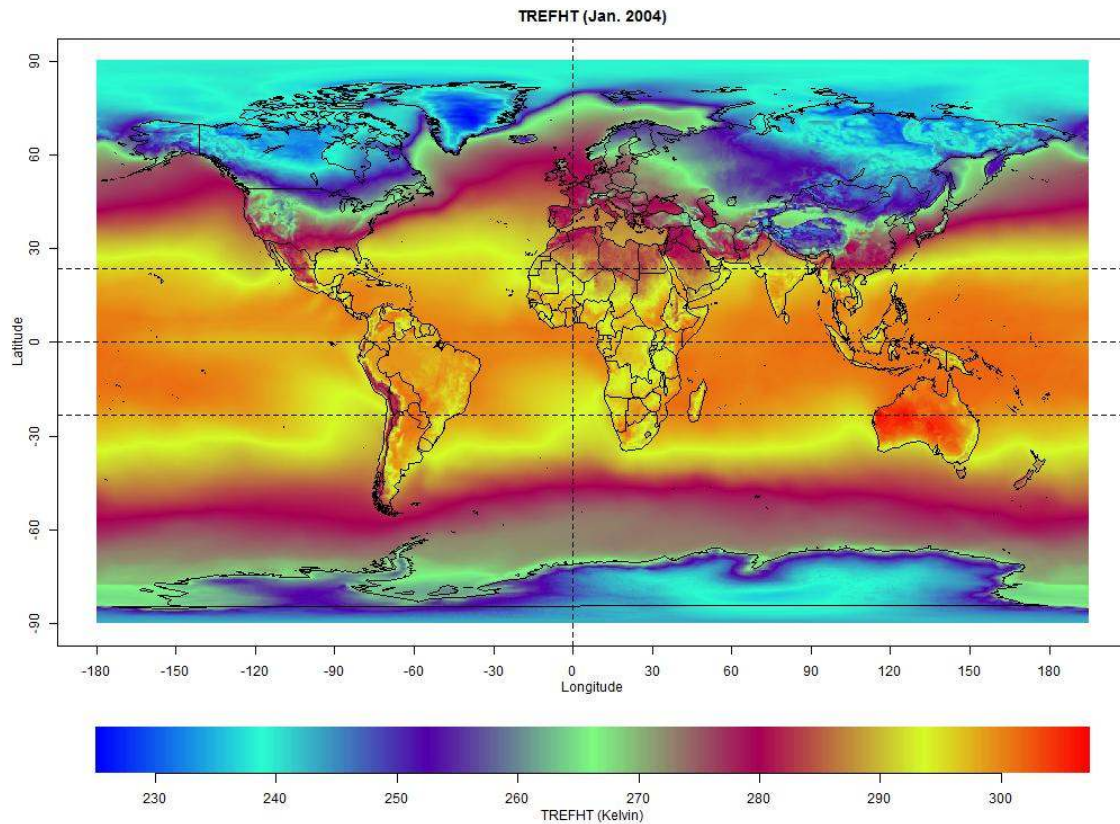


Figure 10.1: Monthly averaged temperature at reference height (TREFHT) in Kelvin (K) for the January 2004. Water freezes at 273.15K and boils at 373.15K.

hemisphere is hotter than the North hemisphere since the seasonal effect.

## 10.2 Parallel Write and Read

*Example: Dump a ddmatrix to a NetCDF4 file and load them from disk.*

The demo command is

Shell Command

```
At the shell prompt, run the demo with 4 processors by
(Use Rscript.exe for windows system)
mpirun -np 4 Rscript -e "demo(nc4_dmat, 'pbdDEMO', ask=F, echo=F)"
```

Main part of the demo is given in the next:

nc4\_dmat

```
1 ### divide data into ddmatrix
2 x <- TREFHT$data
```



```

3 dx <- as.ddmatrix(x)
4
5 # define dimension and variable
6 lon <- ncdim_def("lon", "degree_east", vals =
 TREFHTdefdim[[1]]$vals)
7 lat <- ncdim_def("lat", "degree_north", vals =
 TREFHTdefdim[[2]]$vals)
8 var.def <- ncvar_def("TREFHT", "K", list(lon = lon, lat = lat),
 NULL)
9
10 ### parallel write
11 file.name <- "nc4_dmat.nc"
12 nc <- nc_create_par(file.name, var.def)
13 demo.ncvar_put_dmat(nc, "TREFHT", dx)
14 nc_close(nc)
15 if(comm.rank() == 0){
16 ncdump(file.name)
17 }
18
19 ### parallel read (everyone owns a portion)
20 nc <- nc_open_par(file.name)
21 if(comm.rank() == 0){
22 print(nc)
23 }
24 new.dx <- demo.ncvar_get_dmat(nc, "TREFHT", bldim = bldim(dx),
 ICTXT = dmat.ictxt(dx))
25
26 nc_close(nc)

```

Line 2 and 3 convert `TREFHT$data` into a `ddmatrix` distributed across 4 processors. Line 6 and 7 define the dimensions `lon` and `lat` for longitudes and latitudes, and line 8 defines `var.def` as the dumping variable for “TREFHT” according to the dimensions. Line 12, 13, and 14 create a parallel NetCDF4 file `nc4_dmat.nc`, write the data into the variable on the disk, and close the file. Line 20, 24, and 25 open the file again and read the data from the variable from the data and convert them to a `ddmatrix`.

Note that `demo.ncvar_put_dmat()` and `demo.ncvar_get_dmat()` are implemented for 2D variables only. Please use `pbdNCDF4/ncdf4` primitive functions `ncvar_put()` and `ncvar_get()` via arguments `start` and `count` for more complicated cases. For example, we may write the TREFHT into a slice of a hypercube according to its time step (Jan. 2004).

*Example: Dump and read 9 by 9 (1D and 2D) ddmatrixes in parallel NetCDF4.*

The demo command is

#### Shell Command

```

At the shell prompt, run the demo with 4 processors by
(Use Rscript.exe for windows system)
mpiexec -np 4 Rscript -e "demo(nc4_par_write_1d,'pbdDEMO',ask=F,echo=F)"

```

```
mpiexec -np 4 Rscript -e "demo(nc4_par_write_2d,'pbdDEMO',ask=F,echo=F)"
mpiexec -np 4 Rscript -e "demo(nc4_par_read_1d,'pbdDEMO',ask=F,echo=F)"
mpiexec -np 4 Rscript -e "demo(nc4_par_read_2d,'pbdDEMO',ask=F,echo=F)"
```

These examples create a subdirectory `./nc4_data` first, then generate 81 1D NetCDF4 parallel files `./nc4_data/1d_*.nc`, and 81 2D NetCDF4 parallel files `./nc4_data/2d_*.nc`, and then read them back in. Both reading and writing are in `ddmatrix` and in parallel NetCDF4 formats. Note that the storing formate of nc files are in row major, but R use column major to store a matrix. However, these examples will not have conflicts since they define rows as the frist variable, so transpose matrices will be actually seen from `ncdump`. As long as read and write are in the same order, everything is fine. Further, the `ddmatrix` is actually convert to `gbdr` or `gbdc` formats, then dump to or read from the NetCDF4 files in parallel. This could be inefficient for a large matrix, so use with caution!

### 10.3 Exercises

- 10-1 The demo code `demo/nc4_serial.r` of **pbdDEMO** has a serial version of writing and reading TREHFT as using `ncdf4` on a single NetCDF4 file `nc4_serial.nc`. It is in the sense of single processor programming and has low cost if file is not too big. It is tedious but possible for multiple processors to write a single file with carefully manual barrier and synchronization. Modify `demo/nc4_serial.r` for writing with multiple processors.
- 10-2 It is also possible to read whole chunk of data from a single processor and distribute data later manually. Modify the demo code `demo/nc4_parallel.r` to accomplish this goal and make performance comparisons.
- 10-3 Implement functions or add arguments to the put method, `demo.ncvar_put_dmat()`, and the get method, `demo.ncvar_get_dmat()`, to enable writing and reading high dimension data, for example,  $(lon, lat, time)$  is 2D in time (3D cube) or  $(lon, lat, lev, time)$  is 3D in time (4D hypercube). Dump TREFHT to a slice of 3D cube and load them back to a `ddmatrix`.
- 10-4 In the Sections 11.3 and 11.4, we introduce simple matrix distributed formats `gbdr` and `gbdc` similar to the BLACS contexts ICTXT 2 and 1 with very large block size. The demo code `demo/nc4_gbdc.r` implements similar functionality as for `ddmatrix`, but for `gbdc` format only. Modify the demo code for `gbdr` format. [Hint: See the Exercise 11-4.](#)

*Let no one ignorant of geometry enter here.*

—Plato

One final challenge similar to, but distinct from reading in data is managing data which has already been read into the R processes. Throughout this chapter, we will be making reference to several particulars to the block-cyclic data type used for objects of class `ddmatrix`. In particular, a working knowledge of the block-cyclic data structure and their relationship with BLACS contexts is most useful for the content to follow. As such, the reader is *strongly* encouraged to be familiar with the content of the `pbdDMAT` vignette before proceeding.

## 11.1 Distributed Matrix Redistributions

*Example: Convert between different distributed matrix distributions.*

The demo command is

Shell Command

```
At the shell prompt, run the demo with 4 processors by
(Use Rscript.exe for windows system)
mpiexec -np 4 Rscript -e "demo(reblock,'pbdDEMO',ask=F,echo=F)"
```

The distributed matrix class `ddmatrix` has two components which can be specified, and modified, by the user to drastically affect the composition of the distributed matrix. In particular, these are the object's block-cyclic blocking factor `bldim`, and the BLACS communicator number `CTXT` which sets the 2-dimensional processor grid.

Thankfully, redistributing is a fairly simple process; though we would emphasize that **this is not free of cost**. Reshaping data, especially at scale, can be much more expensive in total than even computation time. That said, sometimes data must move. It is better to get the job done slowly than to “take your ball and go home” with no results. But we caution that if

redistribution can be avoided, then it should, at all costs.

There are several ways one can redistribute a `ddmatrix`. To move the data to a block distribution, one can use the routines `as.rowblock()` and `as.colblock()` for 1-dimensional block distributions, and `as.block()` for a 2-dimensional block distribution. Similarly, there are `as.rowcyclic()` and `as.colcyclic()` functions.

Specifically, these methods take an object of class `ddmatrix` as both an input and an output; i.e., and to emphasize the title of the chapter, this is not a method of *distribution* but *redistribution*. The distribution details of the returned `ddmatrix` are according to the calling method. For example, calling `as.block()` will return a 2-d block-cyclically distributed matrix which is also a 2-d block distributed matrix; see Chapter 5 for information about this distinction.

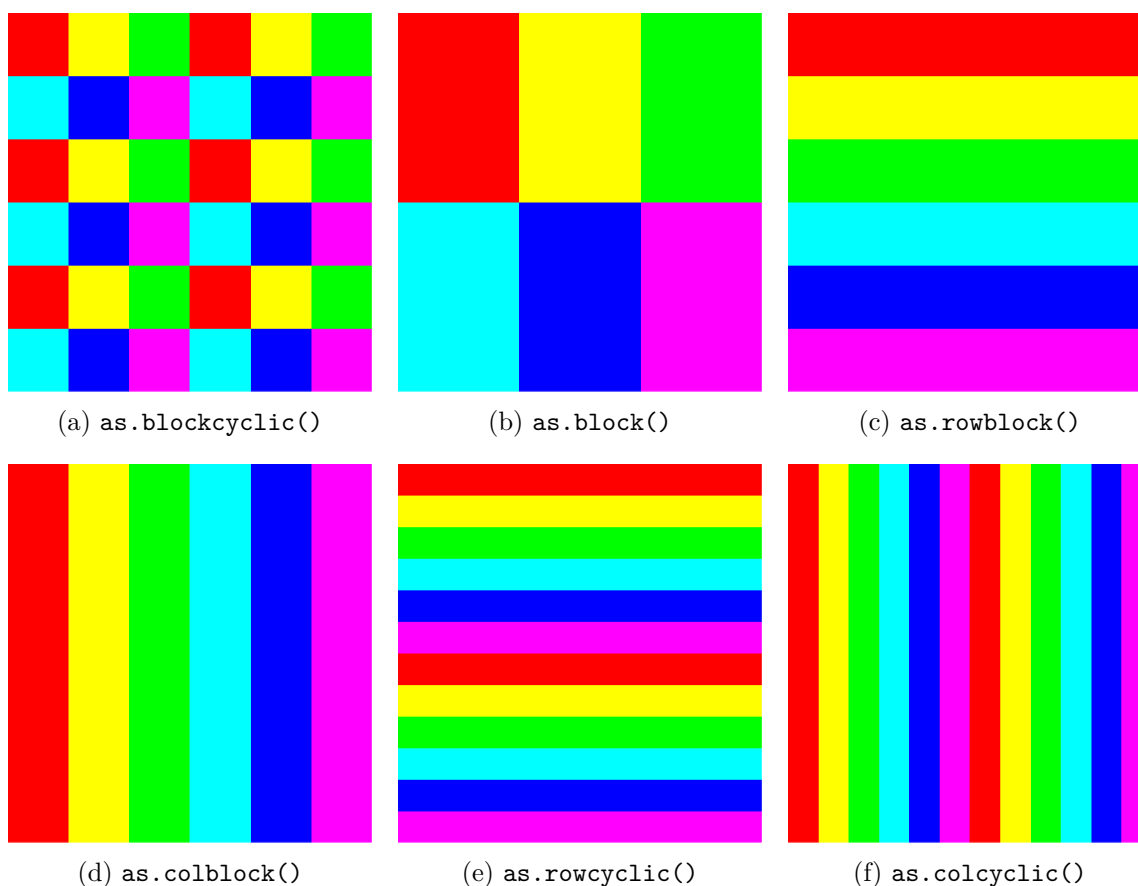


Figure 11.1: Matrix Redistribution Functions

Figure 11.1 shows an example set of outputs for any `ddmatrix` input. Here we assume that there are 6 processors, and in the block and block-cyclic cases, we are assuming that the processor grid (BLACS context) is a  $2 \times 3$  grid.

However, there is a much more general method available, namely `redistribute()`. As the name implies, this method is for reshaping a block-cyclically distributed matrix of one kind to any another. For example, if we have a distributed matrix `dx` and we wish to reshape the distributed matrix so that it now has blocking dimension `newbldim` and is distributed across BLACS context

`newCTXT`, then I need merely call:

```
1 dy <- redistribute(dx, bldim=newbldim, ICTXT=newCTXT)
```

Assuming the data is block cyclic of *any* kind, including degenerate cases, we can convert it to a block cyclic format of any other kind we wish via this `redistribute()` function. The only requirement is that the two different distributions have at least 1 processor in common, and so using the default BLACS contexts (0, 1, and 2) is always acceptable.

## 11.2 Implicit Redistributions

There are several useful functions which apply to distributed matrices, but require a data redistribution as in Section 11, whether the user realizes it or not. These functions are listed in

| Function                  | Example                       | Package        | Effect                               |
|---------------------------|-------------------------------|----------------|--------------------------------------|
| <code>['</code>           | <code>dx[, -1]</code>         | <b>pbdBASE</b> | Row/Column extraction and subsetting |
| <code>na.exclude()</code> | <code>na.exclude(dx)</code>   | <b>pbdBASE</b> | Drop rows with NA's                  |
| <code>apply()</code>      | <code>apply(dx, 2, sd)</code> | <b>pbdDMAT</b> | Applies function to margin           |

Table 11.1: Distributed Matrix Methods with Implicit Data Redistributions

Table 11.1. By default, these functions will re-distribute back to the original data distribution after having performed the initial (necessary) redistribution and performed the requested operations. That is, by default, the problem of managing different data distributions is hidden from the user and entirely implicit. However, there are advantages to becoming familiar with managing these data distributions, because each of these functions has the option to have redistribution directly managed. Now, a data redistribution must occur to use these functions, but understanding which and why can help minimize the number of redistributions performed.

Many of the full details, such as *why* the redistributions need occur in the first place, are outlined in the **pbdDMAT** vignette, but we provide a simple example here. Suppose we have a distributed matrix `dx` distributed on the default grid (i.e., BLACS context 0) and we wish to drop the first column and then use the `apply()` function to extract the p-values, column-wise, of the result of running the Shapiro-Wilk normality test independently on the columns. No one is claiming that this is a wise thing to do, but it is useful for the purpose of demonstration.

To achieve this, we could execute the following:

### Implicit Redistributions

```
1 dx <- dx[-1,]
2
3 result <- apply(dx, MARGIN=2, FUN=function(col)
 shapiro.test(col)$p, reduce=TRUE)
```

In reality, underneath this is actually performing the following sequence of operations:

## Implicit Redistributions

```

1 dx <- redistribute(dx, ICTXT=2)
2 dx <- dx[, -1]
3 dx <- redistribute(dx, ICTXT=0)
4
5 dx <- redistribute(dx, ICTXT=2)
6 result <- apply(dx, MARGIN=2, FUN=function(col)
 shapiro.test(col)$p, reduce=TRUE)

```

Or suppose we wanted instead to drop the first column; then this is equivalent to

## Implicit Redistributions

```

1 dx <- redistribute(dx, ICTXT=1)
2 dx <- dx[, -1]
3 dx <- redistribute(dx, ICTXT=0)
4
5 dx <- redistribute(dx, ICTXT=2)
6 result <- apply(dx, MARGIN=2, FUN=function(col)
 shapiro.test(col)$p, reduce=TRUE)

```

The problem should be obvious. However, thoroughly understanding the problem, we can easily manage the data redistributions using the `ICTXT=` option in these function. So for example, we can minimize the redistributions to only the minimal necessary amount with the following:

## Implicit Redistributions

```

1 dx <- dx[, -1, ICTXT=2]
2
3 result <- apply(dx, MARGIN=2, FUN=function(col)
 shapiro.test(col)$p, reduce=TRUE)

```

This is equivalent to explicitly calling:

## Implicit Redistributions

```

1 dx <- redistribute(dx, ICTXT=2)
2 dx <- dx[, -1, ICTXT=2]
3
4 result <- apply(dx, MARGIN=2, FUN=function(col)
 shapiro.test(col)$p, reduce=TRUE)

```

This is clearly preferred. For more details, see the relevant function documentation.

### 11.3 Load Balance and Unload Balance

*Example: Load balancing (and unbalancing) distributed data.*

The demo command is

#### Shell Command

```
At the shell prompt, run the demo with 4 processors by
(Use Rscript.exe for windows system)
mpiexec -np 4 Rscript -e "demo(balance,'pbdDEMO',ask=F,echo=F)"
```

Suppose we have an unbalanced, distributed input matrix `X.gbd`. We can call `balance.info()` on this object to store some information about how to balance the data load across all processors. This can be useful for tracking data movement, as well as for “unbalancing” later, if we so choose. Next, we call `load.balance()` to obtain a load-balanced object `new.X.gbd`. We can also now undo this entire process and get back to `X.gbd` by calling `unload.balance()` on `new.X.gbd`.

All together, the code looks something like:

#### R Code

```
bal.info <- balance.info(X.gbd)
new.X.gbd <- load.balance(X.gbd, bal.info)
org.X.gbd <- unload.balance(new.X.gbd, bal.info)
```

The details of this exchange are depicted in the example in Figure 11.2. Here, `X.gbd` is unbalanced, and `new.X.gbd` is a balanced version of `X.gbd`.

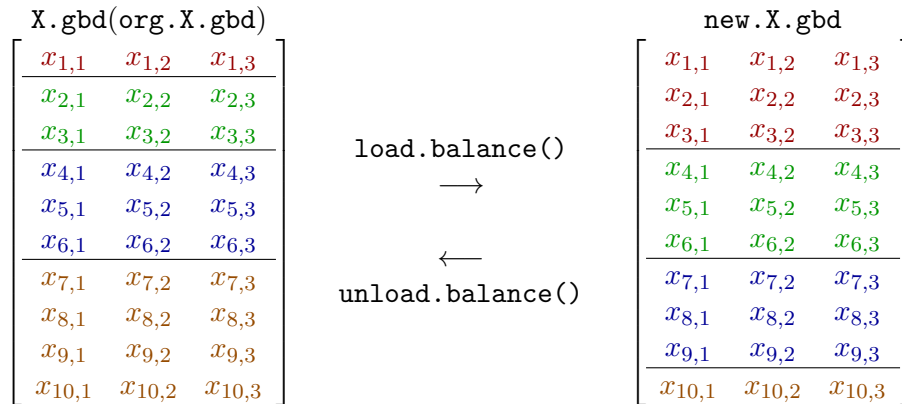


Figure 11.2: Load Balancing/Unbalancing Data:  $\mathbf{X}$  is distributed in `X.gbd(org.X.gbd)` and `new.X.gbd`. Both are distributed row-wise in 4 processors. The colors represent processors 0, 1, 2, and 3, respectively.

The function `balance.info()` is extremely useful, because it will return the information used to load balance the given data `X.gbd`. The return of `balance.info()` is a list consisting of two data frames, `send()` and `recv()`, as well as two vectors, `N.allgbd` and `new.N.allgbd`.

Here, `send` records the original processor rank and the destination processor rank of the unbalanced data (that which is to be transmitted by that processor). The `load.balance()` function uses this table to move the data via `pbdMPI`’s `isend()` function. If any “destination rank” is not the “original rank”, then the corresponding data row will be moved to the appropriate processor. On the other hand, `recv` records the original processor rank and the destination rank of balanced data (that which is received by that processor).

The `N.allgbd` and `new.N.allgbd` objects both have length equal to the communicator containing all numbers of rows of `X.gbd` before and after the balancing, respectively. This is for double checking and avoiding a 0-row matrix issue.

For `unload.balance`, the process amounts to reversing `bal.info` and passing it to `load.balance`.

Finally, note that the “balanced” data is chosen to be balanced in a very particular way; it is arguably not “balanced”, since 3 processors own 3 rows while 1 owns 1 row, and it is perhaps more balanced to have 2 processors own 3 rows and 2 own 2. However, we make this choice for the reason that our “balanced” data will always be a certain kind of degenerate block-cyclic structure. We will discuss this at length in the following section.

## 11.4 Convert Between GBD and DMAT

*Example: Convert between GBD and DMAT formats.*

The demo command is

Shell Command

```
At the shell prompt, run the demo with 4 processors by
(Use Rscript.exe for windows system)
mpirun -np 4 Rscript -e "demo(gbd_dmat,'pbdDMAT',ask=F,echo=F)"
```

The final redistribution challenge we will present is taking an object in GBD format and putting it in the DMAT format. More precisely, we assume the input object `X.gbd` is in GBD and convert the object into an object of class `ddmatrix` which we will call `X.dmat`.

The Figure 11.3 illustrates an example `X.gbd` and `X.dmat` conversion. For full details about the block-cyclic data format used for class `ddmatrix`, see the **pbdDMAT** vignette.

To perform such a redistribution, one simply needs to call:

R Code

```
X.dmat <- gbd2dmat(X.gbd)
```

or

R Code

```
X.gbd <- dmat2gbd(X.dmat)
```

Here, the `gbd2dmat` function does the following:

1. Check number of columns of `X.gbd`. All processors should be the same.
2. Row balance the GBD matrix as necessary via `load.balance()` as in Section 11.3.
3. Call construct a new `ddmatrix` object (via the `new()` constructor) on the balanced matrix, say `X.dmat`, in BLACS context 2 (`ICTXT = 2`).



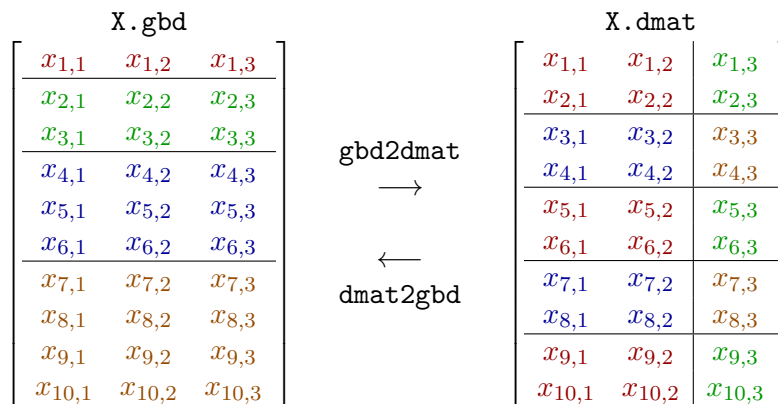


Figure 11.3: Converting Between GBD and DMAT: **X** is distributed in **X.gbd** and **X.dmat**. Both are distributed in 4 processors where colors represents processor 0, 1, 2, and 3. Note that **X.dmat** is in block-cyclic format of  $2 \times 2$  grid with  $2 \times 2$  block dimension.

4. Redistribute **X.dmat** to another BLACS context as needed (default **ICTXT** = 0) via the **base.reblock()** function as in Section 11.1.

Note that the **load.balance()** function, as used above, is legitimately necessary here. Indeed, this function takes a collection of distributed data and converts it into a degenerate block cyclic distribution; namely, this places the data in block “1-cycle” format, distributed across an  $n \times 1$  processor grid. In the context of Figure 11.3 (where the aforementioned process is implicit), this is akin to first moving the data into a distributed matrix format with **bldim=c(3,3)** and **CTXT=2**. Finally, we can take this degenerate block-cyclic distribution and again to Figure 11.3 as our motivating example, we convert the data balanced data so that it has **bldim=c(2,2)** and **CTXT=0**.

## 11.5 Exercises

- 11-1 In the Sections 11.3 and 11.4, we have seen the load balance of GBD matrix and the conversion between GBD and DMAT where GBD matrices **X.gbd** are presumed in row-major as shown in the Figures 11.2 and 11.3. Create new functions **gbdr2gbdc()** and **gbd2gbdr()** converting between row-major and column-major by utilizing functions **gbd2dmat()** and **dmat2gbd()** and changing their option **gbd.major**.
- 11-2 The demo code **demo/gbd\_dmat.r** of **pbdDEMO** has a GBD row-major matrix **X.gbd**. Utilize the functions developed in the Exercise 11-1. Convert **X.gbd** to a column-major matrix **new.X.gbdc** by calling **gbdr2gbdc()**, then convert **new.X.gbdc** back to a row-major matrix **new.X.gbdr** by calling **gbd2gbdr()**. Check if **new.X.gbdr** were the same as **X.gbd**.
- 11-3 In **pbdDEMO**, there are some internal functions **demo.gbdr2dmat()**, **demo.gbd2dmat()**, **demo.dmat2gbdr()**, and **demo.dmat2gbdc()** which have similar implementations as the functions **gbdr2gbdc()** and **gbd2gbdr()** of the Exercise 11-1. Utilize these functions as templates. Create a function **gbd2gbd()** with an argument **new.major** (1, 2) for designated row- or column-majors. Return warnings or errors if the input matrix is not convertible.

- 11-4 The demo code `demo/nc4_gbdc.r` of **pbdDEMO** is an example utilizing GBD column-major matrix `X.gbdc` and dumps the matrix into a NetCDF4 file. Adjust the code. Create a GBD row-major matrix `X.gbdr` and dump the matrix to a new NetCDF4 file `nc4_gbdr.nc` by utilizing the function `ncvar_put_gbd()` with option `gbd.major = 1`. Verify that all `TREFHT` values of both `nc4_gbdc.nc` and `nc4_gbdr.nc` are identical. [Hint: The local matrix of a GBD row- or column-major matrix is still row-major as the default of R.](#)
- 11-5 The `load.balance()` and `unload.balance()` have a potential bug when data size is small and can not fit into the desired block size of a block-cyclic matrix. For instance, four processes in a GBD row-major format with a matrix  $5 \times 1$ . The two functions will (un-)balance the data in  $2 \times 1$  in process 0, and  $1 \times 1$  in others. If the desired block size is 2, then the data should be  $2 \times 1$  in processes 0 and 1,  $1 \times 1$  in process 2, and no element for processor 3. Does any way to fix these two functions?

**Part V**

**Applications**

*Mathematics is the art of giving the same name to different things.*

—Henri Poincare

## 12.1 Introduction

This is a preamble chapter for Chapters 13 and 14, each of which heavily rely on likelihood functions. In a very real sense, likelihoods form the dividing line that separates statistics from other fields, and as such is one of the most important Statistical techniques. The concept of likelihood was popularized in mathematical statistics by R.A. Fisher in 1922 in his landmark paper “On the mathematical foundations of theoretical statistics” (Fisher, 1922). In condensed, broad strokes, likelihood is a tool for developing theoretical inference based on observed data.

We introduce general notations for likelihood functions, which is a standard method for parametric statistics, and is useful for statistical inference (Casella and Berger, 2001). Two useful distributions are introduced. The normal distribution additional to linear model has been applied to the example in Section 4.5. The multivariate normal distribution is also popular to model high dimensional data, and is often used in methods such as model-based clustering in Chapter 13.

Suppose  $\mathbf{X} = \{\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_N\}$  is a random sample, which means the observations are independent and identically distributed (i.i.d.), from a population characterized by a distribution  $\mathcal{F}(\boldsymbol{\theta})$  with unknown parameter  $\boldsymbol{\theta} \in \Theta$ , where  $\Theta$  is the parameter space. Suppose further  $\mathcal{F}$  has a probability density function (pdf for short)  $f(\mathbf{x}_n; \boldsymbol{\theta})$ , with appropriate support. The goal is to estimate  $\boldsymbol{\theta}$  based on the observed data  $\mathbf{x} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$ . Ideally, we want to infer what is the best candidate of  $\boldsymbol{\theta}$  from which we observed  $\mathbf{x}$ . Unlike in Mathematics,  $\mathbf{x}$  is known, but  $\boldsymbol{\theta}$  is unknown and to be determined in Statistics.

A fancy way to estimate  $\boldsymbol{\theta}$  is based on the likelihood function for the observed data  $\mathbf{x}$

$$L(\boldsymbol{\theta}; \mathbf{x}) = \prod_{n=1}^N f(\mathbf{x}_n; \boldsymbol{\theta}) \quad (12.1)$$

or the log likelihood function

$$\log L(\boldsymbol{\theta}; \mathbf{x}) = \sum_{n=1}^N \log f(\mathbf{x}_n; \boldsymbol{\theta}). \quad (12.2)$$

The product on the right hand side of Equation (12.1) is due to the independence assumption of  $\mathbf{X}$ , but the value of  $L(\boldsymbol{\theta}; \mathbf{x})$  may “blow up” to infinity or negative infinity quickly as sample size  $N$  increased. Note that typically, one does not work with the likelihood function in the form of Equation (12.1), but rather the log likelihood of Equation (12.2). The reason for this is that the latter has some nicer properties for most distribution families and is more numerically stable than Equation (12.1).

Statistical methods that deal (directly) with likelihoods involve maximizing (analytically or numerically if the former is not possible or impractical) Equation (12.2) over the parameter space  $\Theta$  to obtain a so-called maximum likelihood estimation, or MLE

$$\hat{\boldsymbol{\theta}}_{ML} := \operatorname{argmax}_{\boldsymbol{\theta} \in \Theta} \log L(\boldsymbol{\theta}; \mathbf{x})$$

Note that the MLE may not exist. There are some additional constraints that are often imposed which make a MLE more well-behaved in some regards, such as regularity conditions of parameter space, or that the parameter  $\boldsymbol{\theta}$  does not depend on the pdf’s support. See [Casella and Berger \(2001\)](#) for details.

## 12.2 Normal Distribution

Section 4.5 offers one way to find  $\boldsymbol{\theta} = \{\boldsymbol{\beta}, \sigma^2\}$  for a linear model without parametric assumption via ordinary least square estimator  $\hat{\boldsymbol{\theta}}_{ols} = \{\hat{\boldsymbol{\beta}}_{ols}, \hat{\sigma}_{ols}^2\}$ . Aside from the Gauss-Markov Theorem, an alternative way is based on likelihood approach by assuming an identical normal distribution with mean zero and variance  $\sigma^2$  to the independent error terms of Equation (4.7). This implies a normal distribution to the response  $y_n$  for  $n = 1, 2, \dots, N$ . More precisely,

$$y_n \stackrel{i.i.d.}{\sim} N(\mathbf{x}_n^\top \boldsymbol{\beta}, \sigma^2) \quad (12.3)$$

where  $\boldsymbol{\theta} = \{\boldsymbol{\beta}, \sigma^2\}$ , and  $\boldsymbol{\beta}$  and  $\mathbf{x}_n$  each have dimension  $p \times 1$ .

By merely mechanically inserting symbols, one may construct a log likelihood function based on the normal density function:

$$\log L(\boldsymbol{\beta}, \sigma^2; \mathbf{y}) = \sum_{n=1}^N \left[ -\frac{1}{2} \log(2\pi\sigma^2) - \frac{(y_n - \mathbf{x}_n^\top \boldsymbol{\beta})^2}{2\sigma^2} \right]. \quad (12.4)$$

The MLEs  $\hat{\boldsymbol{\theta}}_{ML} = \{\hat{\boldsymbol{\beta}}_{ML}, \hat{\sigma}_{ML}^2\}$  can be obtained analytically for this case by taking the first derivatives of Equation (12.4), setting them to zero, and solving the equations. The implementations for numerical solutions (from analytical solutions) or numerical optimization of Equation (12.4) is not difficult and left for the reader in Exercise 12-7.

The assumptions of Statement (12.3) limit the scope of modeling capability, and so next we introduce a more general approach. From the independence assumption and basic multivariate

statistics, statement (12.3) implies a multivariate normal distribution<sup>1</sup> to the response variable  $\mathbf{y}$  with dimension  $N \times 1$ :

$$\mathbf{y} \sim MVN_N(\boldsymbol{\mu}, \boldsymbol{\Sigma}). \quad (12.5)$$

where  $\boldsymbol{\mu} = \mathbf{X}\boldsymbol{\beta}$  with length  $N$ ,  $\boldsymbol{\Sigma} = \sigma^2 \mathbf{I}$  and  $\mathbf{I}$  is an  $N \times N$  identity matrix. In this case,  $\mathbf{y}$  has density function

$$\phi_N(\mathbf{y}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = (2\pi)^{-\frac{N}{2}} |\boldsymbol{\Sigma}|^{-\frac{1}{2}} e^{-\frac{1}{2}(\mathbf{y}-\boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{y}-\boldsymbol{\mu})}$$

and the log likelihood can reduce to Equation (12.4). The MLEs are  $\hat{\boldsymbol{\beta}}_{ML} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$  and  $\sigma_{ML}^2 = \frac{1}{N}(\mathbf{y} - \bar{y}\mathbf{1})^\top (\mathbf{y} - \bar{y}\mathbf{1})$ , where  $\bar{y}$  is the average of  $\mathbf{y}$ , and  $\mathbf{1}$  is the vector of length  $N$  whose entries are all 1.

## 12.3 Likelihood Ratio Test

A very important statistical inference tool based on the likelihood methods discussed so far is the Likelihood Ratio Test (LRT). Provided suitable assumptions hold, this test can be used to compare the fit of two competing models.

Suppose we have data  $\mathbf{X}$  and want to test the hypothesis

$$H_0 : \boldsymbol{\theta} \in \boldsymbol{\Theta}_0$$

against the alternative

$$H_a : \boldsymbol{\theta} \in \boldsymbol{\Theta}_a$$

where the two spaces  $\boldsymbol{\Theta}_0$  and  $\boldsymbol{\Theta}_a$  are not equivalent. The LRT says

$$-2 \log \Lambda(\boldsymbol{\theta}_0, \boldsymbol{\theta}_a; \mathbf{X}) := -2 \log \frac{\max_{\boldsymbol{\theta} \in \boldsymbol{\Theta}_0} L(\boldsymbol{\theta}; \mathbf{X})}{\max_{\boldsymbol{\theta} \in \boldsymbol{\Theta}_0 \cup \boldsymbol{\Theta}_a} L(\boldsymbol{\theta}; \mathbf{X})} \sim \chi_p^2 \quad (12.6)$$

where  $\boldsymbol{\theta}_0$  and  $\boldsymbol{\theta}_a$  are parameters that have maximum likelihoods in spaces  $\boldsymbol{\Theta}_0$  and  $\boldsymbol{\Theta}_0 \cup \boldsymbol{\Theta}_a$  respectively, and  $\chi_p^2$  is a chi-squared distribution with  $p$  degrees of freedom. In some cases,  $p$  is simply the difference in dimension between  $\boldsymbol{\Theta}_0 \cup \boldsymbol{\Theta}_a$  and  $\boldsymbol{\Theta}_0$  (see example below).

For example, in the least squares case of statement (12.5), we may want to test

$$H_0 : \sigma^2 = 1 \quad \text{v.s.} \quad H_a : \sigma^2 > 0$$

which means  $\boldsymbol{\Theta}_0 = \{\boldsymbol{\beta}\}$  and  $\boldsymbol{\Theta}_a = \{\boldsymbol{\beta}, \sigma^2\}$ . Note that  $\boldsymbol{\Theta}_0 \subset \boldsymbol{\Theta}_a$ , so  $\boldsymbol{\Theta}_0 \cup \boldsymbol{\Theta}_a = \boldsymbol{\Theta}_a$ . Given the MLEs  $\hat{\boldsymbol{\theta}}_{0ML}$  and  $\hat{\boldsymbol{\theta}}_{aML}$  for the two spaces  $\boldsymbol{\Theta}_0$  and  $\boldsymbol{\Theta}_a$ , the LRT will be

$$-2 \log \hat{\Lambda}(\hat{\boldsymbol{\theta}}_{0ML}, \hat{\boldsymbol{\theta}}_{aML}; \mathbf{X}) := -2 \log \frac{L(\hat{\boldsymbol{\theta}}_{0ML}; \mathbf{X})}{L(\hat{\boldsymbol{\theta}}_{aML}; \mathbf{X})} \sim \chi_1^2.$$

For type I error  $\alpha = 0.05$ , if the value

$$-2 \log \hat{\Lambda}(\hat{\boldsymbol{\theta}}_{0ML}, \hat{\boldsymbol{\theta}}_{aML}; \mathbf{X}) > q_{\chi_1^2}(0.95) \approx 3.84$$

---

<sup>1</sup>introduced in Section 12.4

where  $q_{\chi^2_1}(0.95)$  is the 95% quantile of chi-squared distribution with 1 degree of freedom. Then we may reject  $H_0 : \sigma^2 = 1$  provided type I error is no greater than 0.05 level.

Note that the LRT introduced here is not dependent on the types of distributions, but has nested parameter space restriction and some regular conditions of parameter space. See [Casella and Berger \(2001\)](#) or [Ferguson \(1996\)](#) for more details of LRTs.

## 12.4 Multivariate Normal Distribution

Suppose  $\{\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_N\}$  is a random sample from multivariate normal distribution (MVN)

$$\mathbf{X}_n \stackrel{i.i.d.}{\sim} MVN_p(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \quad (12.7)$$

where  $\boldsymbol{\theta} = \{\boldsymbol{\mu}, \boldsymbol{\Sigma}\}$ ,  $\boldsymbol{\mu}$  is a center with dimension  $p \times 1$ , and  $\boldsymbol{\Sigma}$  is a  $p \times p$  dispersion matrix. Then  $\mathbf{X}_n$  has density function

$$\phi_p(\mathbf{x}_n; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = (2\pi)^{-\frac{p}{2}} |\boldsymbol{\Sigma}|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\mathbf{x}_n - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x}_n - \boldsymbol{\mu})\right).$$

In general,  $\boldsymbol{\Sigma}$  could be an unstructured dispersion and must be positive definite. Excepting over fitting problems, an unstructured dispersion  $\boldsymbol{\Sigma}$  is desirable to fully characterize correlation of dimensions since the estimation of  $\boldsymbol{\Sigma}$  is completely supported by observed data and there is no restriction on any coordinate of parameter space.

Let  $\mathbf{x} = (\mathbf{x}_1^\top, \mathbf{x}_2^\top, \dots, \mathbf{x}_N^\top)^\top$  be an observed data matrix with dimension  $N \times p$ . The log likelihood function for  $N$  observations is

$$\log L(\boldsymbol{\mu}, \boldsymbol{\Sigma}; \mathbf{x}) = \sum_{n=1}^N -\frac{1}{2} \left[ p \log(2\pi) + \log |\boldsymbol{\Sigma}| + (\mathbf{x}_n - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x}_n - \boldsymbol{\mu}) \right] \quad (12.8)$$

Note that if we wish to numerically compute the log likelihood found in Equation 12.8, the computing time grows as both  $N$  and  $p$  are increased. In some cases, such as model-based clustering in Chapter 13, the total log likelihood is computed in each iteration for all samples and all components.

Suppose  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$  are known. We can efficiently compute the desired quantity using **pbdR**:

R Code

```

1 U <- chol(SIGMA)
2 logdet <- sum(log(abs(diag(U)))) * 2
3 B <- sweep(X.gbd, 2, MU) %*% backsolve(U, diag(1, p))
4
5 # The over- and under-flow need extral care after this step.
6 distval.gbd <- rowSums(B * B)
7
8 distval <- allreduce(sum(distval.gbd))
9 total.logL <- -(p * log(2 * pi) + logdet + distval) * 0.5

```

where  $\mathbf{X}.\text{gbd}$  is a GBD row-major matrix with dimension  $\mathbf{N}.\text{gbd}$  by  $\mathbf{p}$ ,  $\mathbf{MU}$  is a vector of length  $\mathbf{p}$ , and  $\mathbf{SIGMA}$  is a  $\mathbf{p}$  by  $\mathbf{p}$  positive definite matrix. The sample size  $N$  will be the sum of  $\mathbf{N}.\text{gbd}$  across all processors. Note that this trick of computing log likelihood is a one-pass implementation of  $\mathbf{X}.\text{gbd}$ ,  $\mathbf{MU}$ , and  $\mathbf{SIGMA}$ . See HPSC (Chen and Ostrouchov, 2011) or Golub and Van Loan (1996) for more details.

## 12.5 Exercises

12-1 What is the definition of “independent identical distributed”?

12-2 What is the definition of “probability density function”?

12-3 Suppose  $g(\cdot)$  is a continuous function with appropriate support. Argue that  $g(\hat{\theta}_{ML})$  is still a maximum likelihood estimator of  $g(\theta)$ .

12-4 Derive MLEs from Equation (12.4).

12-5 As in Exercise 4-6, argue that  $\hat{\beta}_{ML}$  of Equation (12.4) is also an unbiased estimator of  $\beta$ .

12-6 Show that:

- $\hat{\sigma}_{ML}^2$  of Equation (12.4) is a biased estimator of  $\sigma^2$
- $\hat{\sigma}_{ML}^2$  is an asymptotically unbiased estimator of  $\sigma^2$ .

12-7 Assume data are stored in GBD row-major matrix format. Implement an optimization function for Equation (12.4), numerically optimized via `optim()` in R. Verify the results with the analytical solution.

12-8 Argue that Statement (12.3) implies Statement (12.5), provided appropriated assumption hold.

12-9 Give an example of random variables  $X$  and  $Y$  which are each normally distributed, but  $(X, Y)$  is not a multivariate normal distribution. *Hint:* See Exercise 12-10.

12-10 Show that if  $X$  and  $Y$  independent random variables which are normally distributed, then  $(X, Y)$  has a multivariate normal distribution.

12-11 Prove Statement (12.6). *Hint:* Ferguson (1996).

12-12 Use a similar trick to that of Section 12.4 to implement Principal Component Analysis (PCA). *Hint:* HPSC (Chen and Ostrouchov, 2011).



*If people do not believe that mathematics is simple, it is only because they do not realize how complicated life is.*

—John von Neumann

## 13.1 Introduction

Observational or experimental data are selected or collected with interesting stories, however, deeper discovery enhances values of interpretations. Model-based clustering is an unsupervised learning technique and mainly based on finite mixture models to fit the data, cluster the data, and draw inference from the data (Fraley and Raftery, 2002; Melnykov and Maitra, 2010). The major application of model-based clustering focuses on Gaussian mixture models. For example,  $\mathbf{X}_n$  is a random  $p$ -dimensional observation from the Gaussian mixture model with  $K$  components, which has density

$$f(\mathbf{X}_n; \boldsymbol{\Theta}) = \sum_{k=1}^K \eta_k \phi_p(\mathbf{X}_n; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad (13.1)$$

where  $\phi_p(\cdot; \cdot, \cdot)$  is a  $p$ -dimensional Gaussian/normal density introduced in Section 12.4,

$$\boldsymbol{\Theta} = \{\eta_1, \eta_2, \dots, \eta_{K-1}, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \dots, \boldsymbol{\mu}_K, \boldsymbol{\Sigma}_1, \boldsymbol{\Sigma}_2, \dots, \boldsymbol{\Sigma}_K\},$$

is the parameter space,  $\eta_k$ 's are mixing proportion,  $\boldsymbol{\mu}_k$ 's are the centers of the components, and  $\boldsymbol{\Sigma}_k$ 's are the dispersion of the components.

Suppose a data set  $\mathbf{X} = \{\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_N\}$  has  $N$  observations. Then the log likelihood is

$$\log L(\boldsymbol{\Theta}; \mathbf{X}) = \sum_{n=1}^N \log f(\mathbf{X}_n; \boldsymbol{\Theta}) \quad (13.2)$$

where  $f$  is as in Equation (13.1). Solving the problem of maximizing this log-likelihood is usually done by the expectation-maximization (EM) algorithm (Dempster et al., 1977). Assuming the

EM algorithm converges, we let  $\hat{\Theta}$  be the maximum likelihood estimator of Equation (13.2). Then the maximum posterior probability

$$\operatorname{argmax}_k \frac{\hat{\eta}_k \phi_p(\mathbf{X}_n; \hat{\boldsymbol{\mu}}_k, \hat{\boldsymbol{\Sigma}}_k)}{f(\mathbf{X}_n; \hat{\Theta})}$$

for  $n = 1, 2, \dots, N$  indicates the membership of the observations of the data set  $\mathbf{X}$ .

The **mclust** (Fraley et al., 1999) and **EMCluster** (Chen et al., 2012d) packages are the two main R packages implementing the EM algorithm for the model-based clustering. The **mclust** package has several selections on different kinds of models one may fit, while **EMCluster** implements the most complicated model (dispersions are all unstructured) in a more efficient way, using several initializations, and semi-supervised learning. However, both assume small  $N$  and tiny  $p$ , and only run in serial with sufficient memory.

Note that the k-means algorithm (Forgy, 1965) equivalently assumes  $\eta_1 = \eta_2 = \dots = \eta_K \equiv 1/K$  and  $\boldsymbol{\Sigma}_1 = \boldsymbol{\Sigma}_2 = \dots = \boldsymbol{\Sigma}_K \equiv \mathbf{I}$  in Equation (13.1), where  $\mathbf{I}$  is the identity matrix. As such, the k-means algorithm is a restricted Gaussian mixture model, such that it can be implemented with a simplified version of the EM algorithm. However, due to its strict assumptions, the cluster results are almost always unrealistic, leaving the data scientist unable to draw meaningful inference from the data, and sometimes have unreasonably high classification errors.

## 13.2 Parallel Model-Based Clustering

The **pmclust** (Chen and Ostrouchov, 2012) package is an R package for parallel model-based clustering based on Gaussian mixture models with unstructured dispersions. The package uses data parallelism to solve one very large clustering problem, rather than the embarrassingly parallel problem of fitting many independent models to dataset(s). This approach is especially useful for large, distributed platforms, where the data will be distributed across nodes. And of course it is worth nothing that the package does not merely perform a local clustering operation on the local data pieces; some “gather” and “reduce” operations are necessary at some stages of the parallel EM algorithm.

An expectation-gathering-maximization (EGM) algorithm (Chen et al., 2013) is established for minimizing communication and data movement between nodes. There are four variants of EGM-like algorithms implemented in **pmclust** including EM, AECM (Meng and van Dyk, 1997), APECM (Chen and Maitra, 2011), and APECMa (Chen et al., 2013). The variants are trying to achieve better convergence rates and less computing time than the original EM algorithm. For completeness’ sake, a simple k-means algorithm is also implemented in **pmclust**.

The **pmclust** package is the first **pbdR** application, and the first R package in SPMD to analyze distributed data in Gigabyte scale. It was originally designed for analyzing Climate simulation outputs (CAM5), as discussed in Section 10.1, and is a product for the project “Visual Data Exploration and Analysis of Ultra-large Climate Data” supported by U.S. DOE Office of Science.

The **pmclust** package initially depended on **Rmpi**, but designed in SPMD approach rather than in the manager/worker paradigm even before **pbdR** existed. Later, it migrated to use **pb-**

**dMPI** (Chen et al., 2012a) because of performance issues with **Rmpi** on larger machines. So, by default, the package assumes data are stored in GBD row-major matrix format.

Currently, the package also utilizes **pbdSLAP** (Chen et al., 2012c), **pbdBASE** (Schmidt et al., 2012a), and **pbdDMAT** (Schmidt et al., 2012c) to implement a subset of the above algorithms for data in the **ddmatrix** format. Table 13.1 lists the current implementations.

Table 13.1: Parallel Mode-Based Clustering Algorithms in **pmclust**

| Algorithm                             | GBD | ddmatrix |
|---------------------------------------|-----|----------|
| EM                                    | yes | yes      |
| AECM                                  | yes | no       |
| APECM                                 | yes | no       |
| APECMa                                | yes | no       |
| k-means                               | yes | yes      |
| Based on <b>pmclust</b> version 0.1-4 |     |          |

### 13.3 An Example Using the *Iris* Dataset

The **iris** (Fisher, 1936) dataset is a famous dataset available in R consisting of 50 Iris flowers from each of three species of Iris, namely *Iris setosa*, *Iris versicolor*, and *Iris virginica*. The dataset is tiny, even by today's standards, with only 150 rows and five columns. The column variables consist of the four features sepal length, sepal width, petal length, and petal width, as well as the class of species. We take the first four columns of **iris** to form the matrix  $\mathbf{X}$ , where each row can be classified in three groups by the true id (the fifth column of **iris**) for supervised learning, or clustered in three groups by algorithms for unsupervised learning. Note that the dimension of  $\mathbf{X}$  is  $N = 150$  by  $p = 4$ .

Figure 13.1 shows the pair-wised scatter plot for all features denoted on the diagonal, and classes are indicated by colors. Each panel plots two features on x and y axes. It is clear that **Petal.Length** can split three species in two groups. However, one of the group is mixed with two species and can not be distinguished by any one of these four features.

From the supervised learning point view, the empirical estimation for  $\Theta$  from data will be the best description for the data, assuming the “true model” is a Gaussian mixture. The (serial) demo code **iris\_overlap** in **pbdDEMO** quickly suggests the overlap level of three Iris species. It can be obtained by executing:

R Code

```
R> demo(iris_overlap, 'pbdDEMO', ask = F, echo = F)
```

which utilizes the **overlap** function of **MixSim** (Melnykov et al., 2012). The output is:

R Output

```
R> (ret <- overlap(ETA, MU, S))
$OmegaMap
```

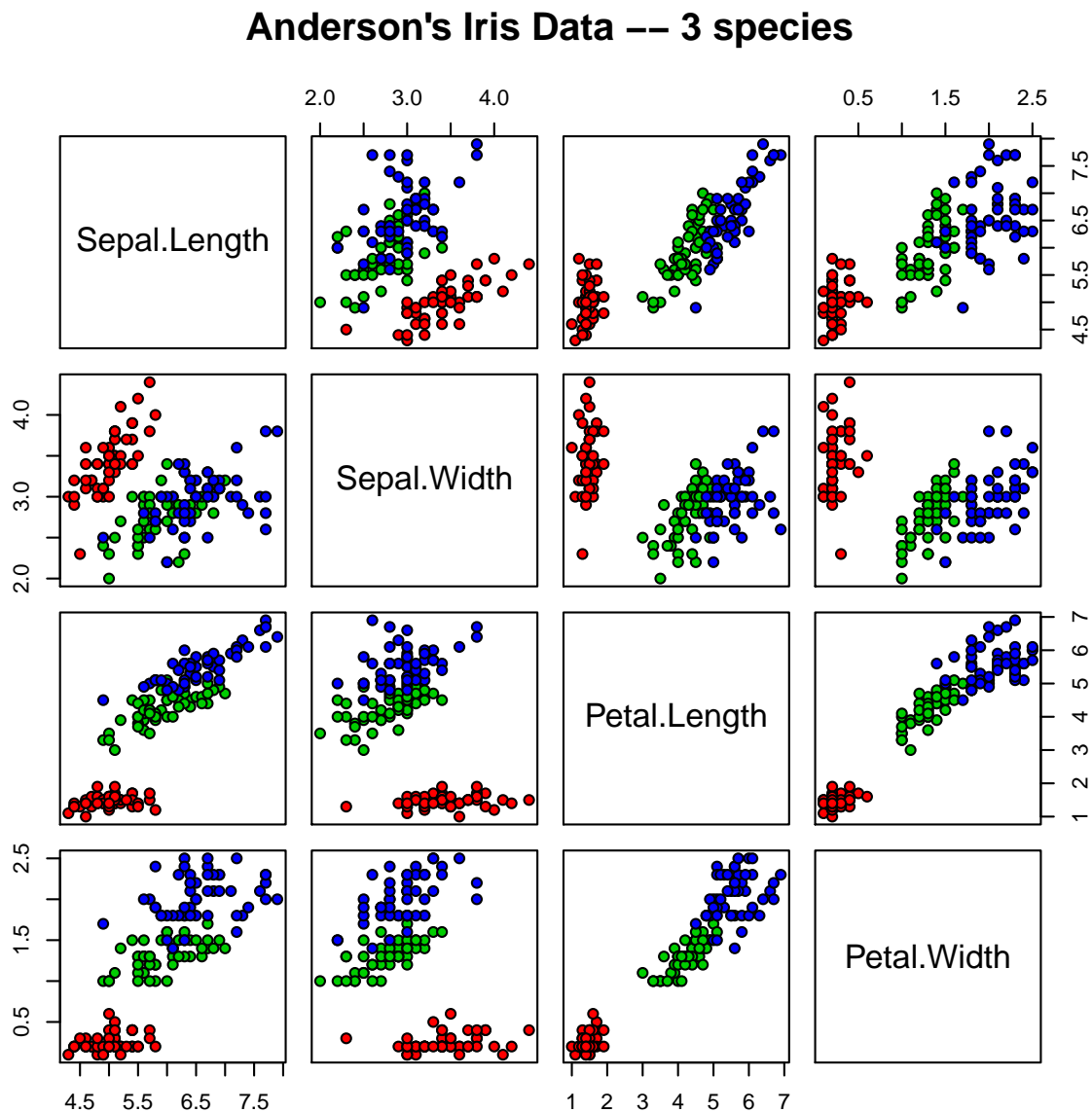


Figure 13.1: Iris pair-wised scatter plot. *Iris setosa* is in red, *Iris versicolor* is in green, and *Iris virginica* is in blue.

```

 [,1] [,2] [,3]
[1,] 1.000000e+00 7.201413e-08 0.00000000
[2,] 1.158418e-07 1.000000e+00 0.02302315
[3,] 0.000000e+00 2.629446e-02 1.00000000

$BarOmega
[1] 0.01643926

$MaxOmega
[1] 0.0493176

```

```

$rcMax
[1] 2 3

R> (levels(iris[, 5]))
[1] "setosa" "versicolor" "virginica"

```

The `OmegaMap` matrix is a map of pair-wise overlap of three species where rows/columns 1, 2, and 3 are *Iris setosa*, *Iris versicolor*, and *Iris virginica*, respectively. The outputs also indicate that the averaged pair-wised overlap (`BarOmega`) is about 1.6%, and the maximum pair-wised overlap (`MaxOmega`) is about 4.9% among these three Iris species. Also, the maximum occurs at 2 (*Iris versicolor*) and 3 (*Iris virginica*) indicating these two species are partly inseparable given these four features.

From the unsupervised learning point view, such as model-based clustering, we must pretend that we are blind to the true class ids, or said another way, we must treat the fifth column of  $\mathbf{X}$  as unobserved. We can then use the four features to form the model and cluster the data, then go back and compare our unsupervised learning results to the true values.

Note that *Iris versicolor* and *Iris virginica* are partly inseparable, so misclassification can happen at the overlap region. We validate the results by comparing the clustering ids to the true class ids using adjusted Rand index (ARI) (Hubert and Arabie, 1985). The ARI takes values between -1 and 1, where 1 is a perfect match. The function `RRand()` in **MixSim** also provides the ARI.

The analysis in the unsupervised learning approach proceeds as follows:

1. decompose  $\mathbf{X}$  on its principal components,
2. project  $\mathbf{X}$  onto the first two principal components (those with largest variability),
3. fit a k-means model and a model-based clustering model, and finally
4. visualize  $\mathbf{X}$  on the plane formed by these new axes, labeling the entries of  $\mathbf{X}$  on this plot with the true ids, and the estimated ids from the clustering algorithms.

This will be the general procedure whether in serial or parallel. For example's sake, we will extend these steps to offer GBD code and `ddmatrix` code to show the similarity of codes.

This example demonstrates that the **pmclust** package can perform analysis correctly, but is not meant to be a demonstration of its scalability prowess. The `iris` dataset is, by any reasonable definition, tiny. Small datasets are generally not worth the trouble of developing parallelized analysis codes for, especially since all the extra overhead costs inherent to parallelism might dominate any theoretical performance gains. Again, the merit of the example is to show off the syntax and accuracy on a single machine; however, **pmclust** scales up nicely to very large dataset running on supercomputers.

### 13.3.1 *Iris* in Serial Code and Sample Outputs

The demo code for the serial Iris example can be found with the package demos, and executed via:

## Shell Command

```
At the shell prompt, run the demo with 4 processors by
(Use Rscript.exe for windows system)
mpirun -np 4 Rscript -e "demo(iris_serial, 'pbdDEMO', ask=F, echo=F)"
```

The code is fairly self-explanatory, and well-commented besides, so we will leave it as an exercise to the reader to read through it carefully.

Running this demo should produce an output that looks something like the following:

[illegible]

```
[1] 0.3333 0.3673 0.2994
null device
 1
```

Finally, Figure 13.2 shows the visualization created by this script.

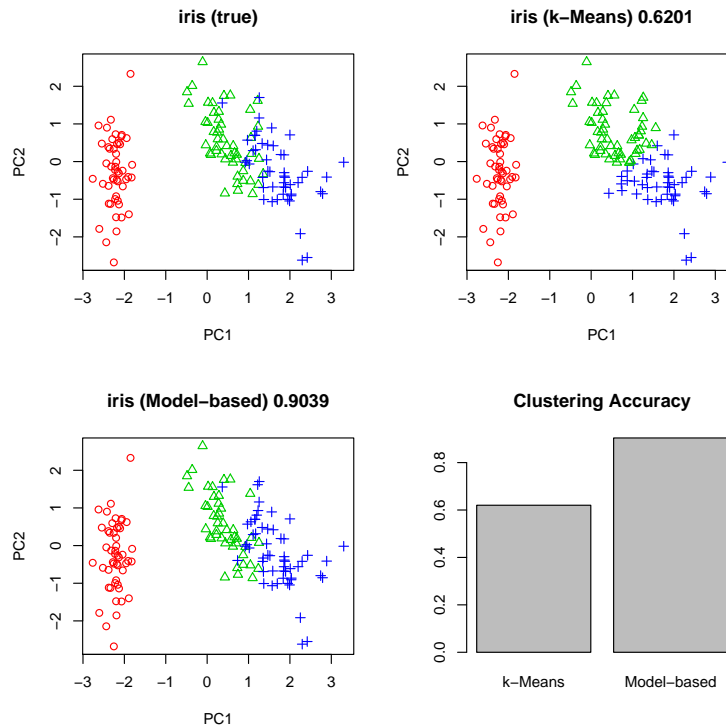


Figure 13.2: Iris Clustering Plots — Serial

### 13.3.2 *Iris* in GBD Code

The demo code for the GBD *Iris* example can be found with the package demos, and executed via:

#### Shell Command

```
At the shell prompt, run the demo with 4 processors by
(Use Rscript.exe for windows system)
mpirun -np 4 Rscript -e "demo(iris_gbd, 'pbdDEMO', ask=F, echo=F)"
```

### Sample Outputs

Running this script should produce an output that looks something like the following:

```

COMM.RANK = 0
[1] 2.547376e-14 8.076873e-15 4.440892e-14
COMM.RANK = 0
[1] 0.6201352 0.6311581 0.6928082
null device
 1

```

Finally, figure 13.3 shows the visualization created by this script.

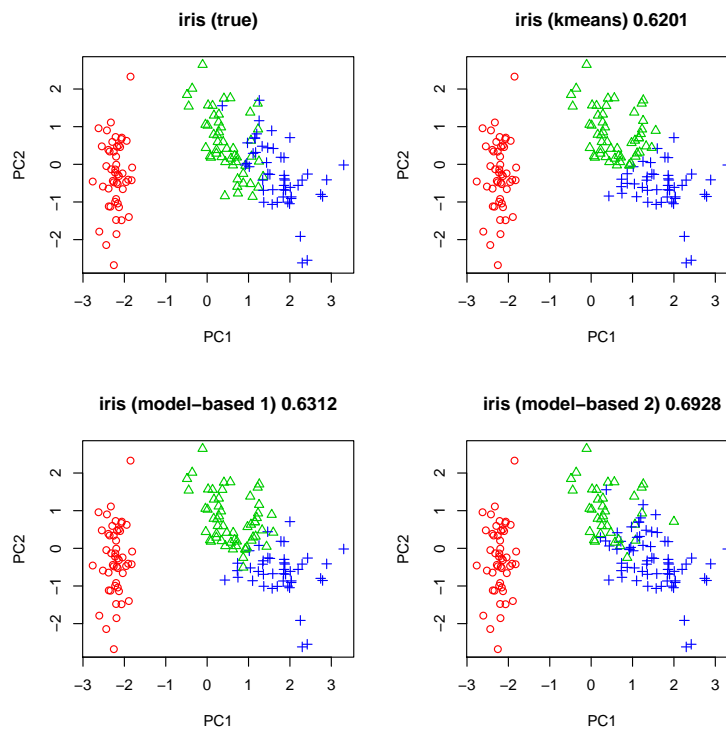


Figure 13.3: Iris Clustering Plots — GBD

### 13.3.3 *Iris* in ddmatrix Code

The demo code for the DMAT Iris example can be found with the package demos, and executed via:

Shell Command

```

At the shell prompt, run the demo with 4 processors by
(Use Rscript.exe for windows system)
mpirun -np 4 Rscript -e "demo(iris_dmat, 'pbdDEMO', ask=F, echo=F)"

```



## Sample Outputs

Running this script should produce an output that looks something like the following:

```
Using 2x2 for the default grid size

COMM.RANK = 0
 [,1] [,2] [,3] [,4]
[1,] -4.440892e-16 1.990595e-16 -2.428613e-17 2.498002e-16
COMM.RANK = 0
 [,1] [,2] [,3] [,4]
[1,] 1.0000000 -0.1175698 0.8717538 0.8179411
[2,] -0.1175698 1.0000000 -0.4284401 -0.3661259
[3,] 0.8717538 -0.4284401 1.0000000 0.9628654
[4,] 0.8179411 -0.3661259 0.9628654 1.0000000
COMM.RANK = 0
[1] 0.645147

null device
 1
```

Finally, figure 13.3 shows the visualization created by this script.

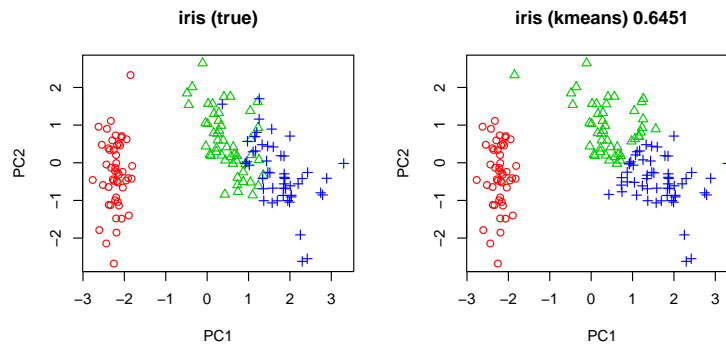


Figure 13.4: Iris Clustering Plots — GBD

## 13.4 Exercises

13-1 As Figures 13.2 and 13.3, none of clustering method is able to obtain the true. However, there are ways that may improve the final clustering and close to the true, including

- 1) by reducing the convergent criteria,
- 2) by increasing the number of initialization steps,
- 3) by aggregating several initialization strategies, and
- 4) by given some prior information about classification.

Using `iris` as a data, and trying different ways to see if final clustering results are improved. See the next Exercises for details.

13-2 In serial, utilizing **MixSim** to generate parameters with different levels of overlaps, based on the parameters to generated data from Gaussian mixture models, and repeat Exercise 13-1 on the generated data to show how overlaps can affect algorithm performances by comparing ARIs.

13-3 In serial, utilizing **EMCluster** on the generated data from Exercise 13-2 to show and test how initialization strategies can affect algorithm performances by comparing ARIs.

13-4 In serial, **EMCluster** also implements semi-supervised model-based clustering, select some high density points from the generated data and labeling them as prior know information, then test how these information can affect algorithm performances by comparing ARIs.

13-5 Argue that the k-means algorithm (Forgy, 1965) equivalently assumes  $\eta_1 = \eta_2 = \dots = \eta_K \equiv 1/K$  and  $\Sigma_1 = \Sigma_2 = \dots = \Sigma_K \equiv \mathbf{I}$  in Equation (13.1), where  $\mathbf{I}$  is the identity matrix.

13-6 The EM algorithm is a typical way to find the MLEs for mixture models. State the two steps (E- and M-steps) of the EM algorithm in general, and argue the monotonicity of log likelihood in every iteration. [Hint: Jensen's Inequality \(Jensen, 1906\).](#)

## Phylogenetic Clustering (Phyloclustering)

*The scientific imagination always restrains itself  
within the limits of probability.*

—Thomas Huxley

### 14.1 Introduction

Phylogenetic Clustering (Phyloclustering) discovers population structure based on information of DNA/RNA sequences by combining two inventions: model-based clustering with evolutionary models (Chen, 2011). Note that what speaking here, regarding to “evolutionary”, is a mathematical/statistical model to interpret biological targets. Neither religion nor theology is involved.

In an over simplified case, suppose a sequence is composed by four nucleotides  $\mathcal{S} = \{\text{A}, \text{G}, \text{C}, \text{T}\}$ . Assume a sequence  $\mathbf{x}_n = \{x_{n1}, x_{n2}, \dots, x_{nL}\} \in \mathcal{S}$  has  $L$  loci (positions ordered) and is observed from a population, but may have  $K$  subpopulations that similar sequence patterns are expected within each common subpopulation. Each subpopulation is represented by a common center sequence  $\boldsymbol{\mu}_k = \{\mu_{k1}, \mu_{k2}, \dots, \mu_{kL}\} \in \mathcal{S}$  which may or may not hypothetically exit in population and has to be determined. Therefore, each sequence has a probability mutated/evolved from any center sequence. The higher the probability, the closer (more similar) to the center sequence. This bold assumption may be invalid to and even violate traditional phylogeny construction and evolutionary research, but it is a comparative way to reconstruct population structures totally based on the discovered facts of observed data.

The evolutionary model is based on a continuous time Markov chain (CTMC) model on a state space  $\mathcal{S}$  that the mutation process is characterized by an instantaneously rate matrix  $\mathbf{Q}$  with dimension  $4 \times 4$ , i.e. rate at scale of tiny mutation time  $t \rightarrow 0$ . We use the following steps to construct the likelihood function as introduced in Chapter 12:

1. Given the above setting, the mutation chance from a nucleotide  $x$  to a nucleotide  $y$  in time  $t$  is

$$\mathbb{P}_{x,y}(t) = e^{\mathbf{Q}_{x,y}t} \quad (14.1)$$

for all  $x, y \in \mathcal{S}$ .

2. Assume each locus is mutated independently, then the mutation chance (the transition probability) from  $\mu_k$  to  $x_n$  in time  $t$  is

$$p_{\mu_k, x_n}(t) = \prod_{l=1}^L \mathbb{P}_{\mu_{kl}, x_{nl}}(t)$$

for all  $\mu_{kl}, x_{nl} \in \mathcal{S}$ .

3. Suppose there are  $K$  subpopulations with mixing proportion  $\eta_k$ 's, then the mutation chance from a sequence  $\mu_k$  to a sequence  $x_n$  is

$$f(x_n; \theta_K) = \sum_{k=1}^K \eta_k p_{\mu_k, x_n}(t) \quad (14.2)$$

where  $\theta_K = \{\eta_1, \eta_2, \dots, \eta_{K-1}, \mu_1, \mu_2, \dots, \mu_K, Q, t\}$  are unknown and to be determined. For simplicity, assume  $Q$  and  $t$  are identical across  $K$  subpopulations. Denote the distribution  $\mathcal{F}(\theta_K)$  of the density function  $f(x_n; \theta_K)$  for  $x_n$ .

4. Suppose observed  $N$  sequences  $\mathbf{x} = \{x_1, x_2, \dots, x_N\}$  (each has  $L$  loci) independently and identically selected from unknown  $K$  subpopulations with mixing proportion  $\eta$  to be estimated, then the likelihood is

$$L(\theta_K; \mathbf{x}) = \prod_{n=1}^N f(x_n; \theta_K).$$

See Section 12.1 for construction.

5. In short, the log likelihood is

$$\begin{aligned} \log L(\theta_K; \mathbf{x}) &= \sum_{k=1}^K \log f(x_n; \theta_K) \\ &= \sum_{k=1}^K \log \left[ \sum_{k=1}^K \eta_k p_{\mu_k, x_n}(t) \right] \\ &= \sum_{k=1}^K \log \left[ \sum_{k=1}^K \eta_k \left( \prod_{l=1}^L \mathbb{P}_{\mu_{kl}, x_{nl}}(t) \right) \right] \\ &= \sum_{k=1}^K \log \left[ \sum_{k=1}^K \eta_k \left( \prod_{l=1}^L e^{Q_{\mu_{kl}, x_{nl}} t} \right) \right]. \end{aligned} \quad (14.3)$$

Equation (14.2) has similar structure as Equation (13.1). Therefore, the EM algorithm (Dempster et al., 1977) can be applied to maximize Equation (14.3) as maximize Equation (13.2). Except the parameter space  $\Theta_K$  of Equation (14.3) where  $\theta_K$  belongs to is neither continuous nor discrete space since  $x_n$  and  $\mu_k$  are in a categorical space which yields a very different E- and M-steps.

## 14.2 The phyclus Package

The **phyclus** (Chen, 2011) is an R package fully implements phyloclustering with different configurations, EM algorithms, and incorporating several useful tools such as **ms** (Hudson, 2002) for simulating phylogeny and **seq-gen** (Rambaut and Grassly, 1997) for simulating sequence with vary mutations based on phylogenies. The **phyclus** also provides functions for re-sampling sequences from predicted models for determining an appropriate number of subpopulations. Those functions are particular useful for Sections 14.3 and 14.4.

The **phyclus** package has several example datasets which is initialed by several longitudinal animal studies on Equine Infectious Anemia Virus (EIAV) (Leroux et al., 2004). The EIAV is a lentivirus that infects equine and causes Equine Infectious Anemia (EIA), and it is similar to Human Immunodeficiency Virus (HIV) infects human and causes Acquired Immunodeficiency Syndrome (AIDS). Figure 14.1 (Weiss, 2006) shows a phylogeny of several relative lentivirus in the retrovirus family, it also shows the closeness of EIAV and HIV which makes the possible to build an animal model based on EIAV and to study viral transmission mechanism further in HIV.

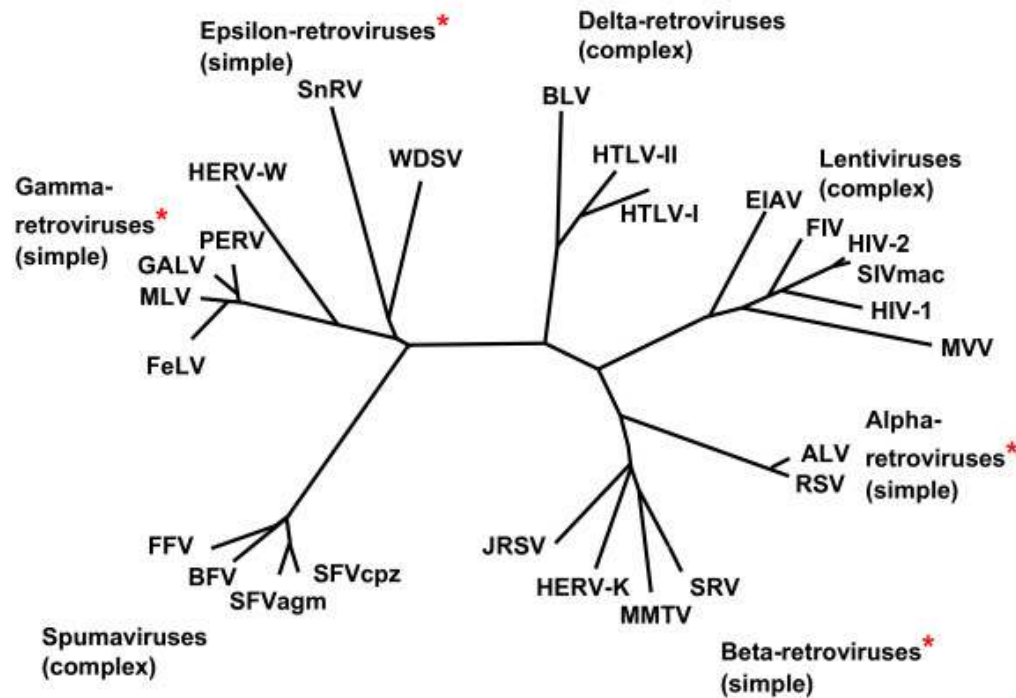


Figure 14.1: Retrovirus phylogeny originated from Weiss (2006).

The disease EIA progresses as the immune system response to the viruses population change in blood which is collected over time and generations. Part of blood samples associated with fever cycles are sequenced to identify highly mutable coding regions with several overlapping reading frames. Immune system response to new mutants of EIAV and trigger fever as a major

signal and symptom of EIA. Therefore, the sequences and regions then can be associated with disease progresses for further analysis. Identify population structures is the critical step for understanding the mutation patterns and designing better medicine or vaccine.

We perform phyloclustering on an example dataset, *Pony 524* (Carpenter et al., 2011), which is given in Figure 14.2. See Baccam et al. (2003) for more about the studies and stories of infected horses. It plots the example dataset where  $N = 146$  EIAV sequences are in y-axis and  $L = 405$  loci in x-axis. The top row is the consensus sequence, and only mutation sites are spotted for 146 sequences. Colors represent A, C, G, and T nucleotides. Three clusters fitted by a CTMC model are shown where common mutation locations and types are grouped by colored spots.

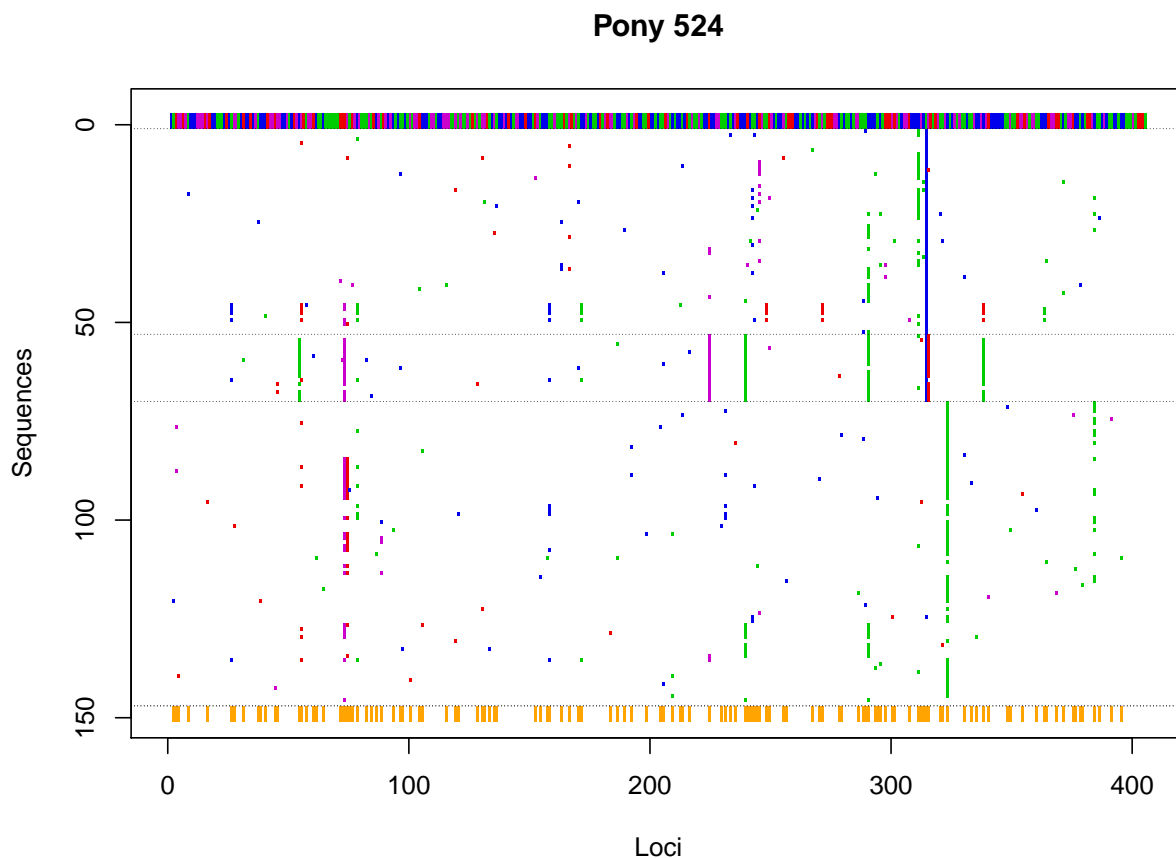


Figure 14.2: 146 EIAV sequences of *Pony 524* in three clusters.

### 14.3 Bootstrap Method

“How many clusters are appropriate for the data?” is a typical question to any good scientists. There are several ways trying to infer this from data in statistics via hypothesis testing. For example,  $H_0 : K = 2$  v.s.  $H_a : K = 3$  or more generally  $H_0 : K = K'$  v.s.  $H_a : K = K^*$  for

any  $K' < K^*$ . In mixture models, the nested parameter space is inappropriate, hence, the LRT introduced in Section 12.3 may not appropriate.

The Bootstrap method (Efron, 1979) may provide an adequate solution to rebuild an asymptotic distribution for the likelihood ratio (LR). “Asymptotic” means ideally large sample size property. The bootstrap method is a re-sampling technique based on Monte Carlo property either from data (non-parametric) or from model (parametric) to form a distribution for a testing statistics. Here we need a distribution such as a hypothetically chi-squared distribution for LR where the degrees of freedom are difficult to be determined. Therefore, we may obtain a p-value by comparing LR to this distribution rather than deriving an asymptotic distribution from LRT.

Phyloclustering which uses a mixture models with unusual parameter space which is also particular suitable to apply the bootstrap methods to determine an appropriate number of subpopulations. For given data  $\mathbf{X}$  and hypothetical  $K'$  and  $K^*$ , we may perform parametric bootstrap as the next.

**Step 1:** Based on  $\mathbf{X}$ , obtain MLEs  $\hat{\boldsymbol{\theta}}_{K' ML}$  and  $\hat{\boldsymbol{\theta}}_{K^* ML}$  under  $\boldsymbol{\Theta}_{K'}$  and  $\boldsymbol{\Theta}_{K^*}$ , respectively.

**Step 2:** Compute and let  $\hat{\lambda} := -2 \log \hat{\Lambda}(\hat{\boldsymbol{\theta}}_{K' ML}, \hat{\boldsymbol{\theta}}_{K^* ML}; \mathbf{X})$ .

**Step 3:** Sample new data  $\mathbf{X}^{(b)}$  from  $\mathcal{F}(\hat{\boldsymbol{\theta}}_{K' ML})$ .

**Step 4:** Based on  $\mathbf{X}^{(b)}$ , obtain MLEs  $\hat{\boldsymbol{\theta}}_{K' ML}^{(b)}$  and  $\hat{\boldsymbol{\theta}}_{K^* ML}^{(b)}$  under  $\boldsymbol{\Theta}_{K'}$  and  $\boldsymbol{\Theta}_{K^*}$ , respectively, via the EM algorithm.

**Step 5:** Compute and let  $\lambda^{(b)} := -2 \log \hat{\Lambda}(\hat{\boldsymbol{\theta}}_{K' ML}^{(b)}, \hat{\boldsymbol{\theta}}_{K^* ML}^{(b)}; \mathbf{X}^{(b)})$ .

**Step 6:** Repeat Steps 3 to 5 for  $B$  times, collect and let  $\mathcal{F}^{(B)}(\lambda) := \{\lambda^{(1)}, \lambda^{(2)}, \dots, \lambda^{(B)}\}$  which is an approximation to  $\mathcal{F}(\lambda)$ , the distribution of  $\lambda$ , as  $B$  large enough.

**Step 7:** If  $\hat{\lambda}$  is greater than  $q_{\mathcal{F}^{(B)}(\lambda)}(0.95)$ , then we reject the  $K'$  model under 0.05 level of type I error with  $B$  bootstrap samples.

Unlike LRT of Section 12.3, note that  $\hat{\boldsymbol{\theta}}_{K^* ML}$  in Step 1 and  $\hat{\boldsymbol{\theta}}_{K^* ML}^{(b)}$  in Step 4 are MLEs in space  $\boldsymbol{\Theta}_{K^*}$  rather than the spaces  $\boldsymbol{\Theta}_{K'} \cup \boldsymbol{\Theta}_{K^*}$  nor  $\boldsymbol{\Theta}_{K'+K^*}$ , which means no guarantee the estimators are the MLEs of larger spaces. This makes the general LRT invalid for mixture models, therefore, other information criteria such as AIC (Akaike, 1974) are also questionable for determining a suitable number of clusters. Parametric or non-parametric bootstraps are other robust methods to verify and provide a suggestion. See Chen (2011) for more simulation studies of this approach via **phyclust**.

## 14.4 Task Pull Parallelism

Obviously, Step 4 will be computationally intensive as  $B$  increased, and no guarantee that each of  $b = 1, 2, \dots, B$  bootstrap samples will take similar time at obtaining MLEs. It may be possible to parallelize the EM algorithm fully in SPMD such as Section 13.2, however, this step is still a bottleneck of whole computation in general.

The task parallelism as mention in Exercise 2-2 is one way to solve the problem by simply divided jobs equally likely to all processors. This is probably an optimal solution for equal loading jobs in homogeneous computing environment. However, it will be a terrible solution for unbalance loading jobs or in-homogeneous computing environment, such as bootstrap methods introduced in Section 14.3. Note that there are also some drawbacks for task parallelism:

- it requires one processor to handle job controls as the role of manager in manager/workers programming paradigm,
- the code is not obviously and difficult to debug or generalize,
- the code requires further reordering for returned results, and
- jobs may break in workers which can cause crash of entire computation.

The website at <http://math.acadiau.ca/ACMMaC/Rmpi/examples.html> has a general view of task parallelism and examples in **Rmpi**. Among three task parallel methods, task pull has the best performance and suit for bootstrap methods. A simplified example of task pull in SPMD can be found in the **pbdMPI** demo via

```
At the shell prompt, run the demo with 4 processors by
(Use Rscript.exe for windows system)
mpiexec -np 4 Rscript -e "demo(task_pull,'pbdMPI',ask=F,echo=F)"
```

which does the following

#### Task Pull Example in **pbdMPI**

```
1 ### Initial
2 library(pbdMPI, quiet = TRUE)
3
4 ### Examples
5 FUN <- function(jid){
6 Sys.sleep(1)
7 jid * 10
8 }
9
10 ret <- task.pull(1:10, FUN)
11 comm.print(ret)
12
13 if(comm.rank() == 0){
14 ret.jobs <- unlist(ret)
15 ret.jobs <- ret.jobs[names(ret.jobs) == "ret"]
16 print(ret.jobs)
17 }
18
19 ### Finish
20 finalize()
```

Lines 5 to 8 define a major function to be evaluated on all workers which are ranks 1, 2, and 3 in



this case. Line 10 prepares 10 jobs from 1 to 10 where jobs can be done by any available worker. The `task.pull()` is actually a combination of two functions `task.pull.workers()` called by all workers and `task.pull.master()` only called by the master by default rank 0. Lines 13 to 17 extract and summarize all returned results on master.

## 14.5 An Example Using the *Pony 524* Dataset

As introduced in Section 14.2, we will fit the  $K' = 1$  and  $K^* = 2$  first. Then, we use bootstrap method in Section 14.3 to find out better number of clusters based on  $B = 100$ , i.e. we compare  $\hat{\lambda} = -2\log \hat{\Lambda}(\hat{\theta}_{1,ML}, \hat{\theta}_{2,ML}; \mathbf{X})$  to  $\lambda^{(b)} = -2\log \hat{\Lambda}(\hat{\theta}_{1,ML}^{(b)}, \hat{\theta}_{2,ML}^{(b)}; \mathbf{X}^{(b)})$  with  $b = 1, 2, \dots, B$  bootstrap samples  $\mathbf{X}^{(b)}$  which are generated from  $\hat{\theta}_{1,ML}$ . The idea here is if one cluster (or non cluster) were suggested from data, then  $\hat{\lambda}$  would have tiny chance located at the tail region of  $\lambda^{(b)}$ 's histogram. On the other hand, if  $\hat{\lambda}$  did located at the tail region, then we might say it is not happen randomly and the evidence is significant to reject one cluster with small error when comparing with two clusters based on  $B$  bootstrap samples.

The whole processes are designed using the task pull method in Section 14.4 to efficiently complete all likelihood estimations of all bootstrap cases on 4 processors (1 master and 3 workers). The demo code for the Pony 524 example can be found with the package demos, and executed via:

```
At the shell prompt, run the demo with 4 processors by
(Use Rscript.exe for windows system)
mpiexec -np 4 Rscript -e
 "demo(phyclust_bootstrap, 'pbdDEMO', ask=F, echo=F)"
```

After a long running, this demo should produce an output that looks something like the following:

### R Output

```
K0: 1
Ka: 2
logL K0: -4033.154
logL Ka: -3403.873
LRT: -1258.561
p-value: 0
```

Note that rerun the code again may produce different results since job assignments to workers are randomly dependent on bootstrap samples. However, the conclusion should be similar that  $K' = 1$  is rejected under 0.05 level type I error using  $B = 100$  parametric bootstraps samples. This may suggest more than one subpopulations exist in this dataset, and more detail investigation should be conducted.

## 14.6 Exercises

- 14-1 Argue that the instantaneous rate matrix  $\mathbf{Q}$  of Equation 14.1 is positive definite. Therefore, argue that the eigenvalue decomposition of  $\mathbf{Q} = \mathbf{U}\mathbf{D}\mathbf{U}^{-1}$  exists. Prove that  $e^{\mathbf{Q}t} = \mathbf{U}e^{\mathbf{D}t}\mathbf{U}^{-1}$ . Hence, this is an easy way for computing transition probability  $\mathbb{P}_{x,y}(t)$ .
- 14-2 Argue that  $\mathcal{F}^{(B)}(\lambda)$  is a good approximation to  $\mathcal{F}(\lambda)$  in Step 4 of Section 14.3.
- 14-3 In Section 14.5, we have tested  $H_0 : K' = 1$  v.s.  $H_a : K^* = 2$ . Change the code to test  $H_0 : K' = 1$  v.s.  $H_a : K^* = 3$  and  $H_0 : K' = 2$  v.s.  $H_a : K^* = 3$ . Draw conclusions for these tests.
- 14-4 Implement a function `task.push()` for task push parallelism utilizing examples at the website <http://math.acadiau.ca/ACMMaC/Rmpi/examples.html>.
- 14-5 Compare the computation time of `parLapply()`, `task.pull()`, and `task.push()` using in SPMD and in 4 processors. By testing  $H_0 : K' = 1$  v.s.  $H_a : K^* = 2$  to Pony 524 dataset using bootstrap method as Section 14.5.

*A good Bayesian does better than a non-Bayesian, but a bad Bayesian gets clobbered.*

—Herman Rubin

## 15.1 Introduction

In modern statistics, likelihood principle introduced in Chapter 12 has produced several advantages to data analysis and statistical modeling. However, as model getting larger and data size getting bigger, the maximization of likelihood function becomes infeasible analytically and numerically. Bayesian statistics based on Bayes theorem somehow relieves the burden of optimization, but it changes the way of statistical inference.

In likelihood principle, we based on maximum likelihood estimators for estimations, hypothesis testings, confidence intervals, etc. In Bayesian framework, we make inference based on posterior distribution, which is a composition of likelihood and prior information, such as for posterior means and credible intervals. For more information about Bayesian statistics, readers are encouraged to read Berger (1993); Gelman et al. (2003).

Mathematically, we denote  $\pi(\boldsymbol{\theta}|\mathbf{x})$  for posterior,  $p(\mathbf{x}|\boldsymbol{\theta})$  for likelihood, and  $\pi(\boldsymbol{\theta})$  for prior where  $\mathbf{x}$  is a collection of data and  $\boldsymbol{\theta}$  is a set of interesting parameters. The idea of Bayes theorem says

$$\pi(\boldsymbol{\theta}|\mathbf{x}) = \frac{p(\mathbf{x}|\boldsymbol{\theta})\pi(\boldsymbol{\theta})}{\int p(\mathbf{x}|\boldsymbol{\theta})\pi(\boldsymbol{\theta})d\boldsymbol{\theta}} \quad (15.1)$$

$$\propto p(\mathbf{x}|\boldsymbol{\theta})\pi(\boldsymbol{\theta}) \quad (15.2)$$

in short, the posterior is proportional to the product of likelihood and prior. Note that the integral denominator of Equation (15.1) can be seen as a normalizing constant, and is usually ignorable in most of Bayesian calculation, then Equation (15.2) provides great reduction tricks for analytical and simulated solutions.

For example, suppose  $\mathbf{x} = \{x_1, x_2, \dots, x_n\}$  are random samples from  $N(\mu, \sigma^2)$  where  $\mu$  is unknown and needed to be inferred (i.e.  $\boldsymbol{\theta} = \{\mu\}$ ), and  $\sigma^2$  is known. Suppose further  $\mu$  has a prior

distribution  $N(\mu_0, \sigma_0^2)$  where  $\mu_0$  and  $\sigma_0^2$  are hypothetically known. After a few calculation, we have the posterior for  $\mu|\mathbf{x}$  denoted by conventional syntaxes next.

$$\mathbf{x} \stackrel{i.i.d.}{\sim} N(\mu, \sigma^2) \quad (15.3)$$

$$\mu \sim N(\mu_0, \sigma_0^2) \quad (15.4)$$

$$\mu|\mathbf{x} \sim N(\mu_n, \sigma_n^2) \quad (15.5)$$

where  $\mu_n = \sigma_n^2 \left( \frac{\mu_0}{\sigma_0^2} + \frac{n\bar{x}}{\sigma^2} \right)$ ,  $\sigma_n^2 = \left( \frac{1}{\sigma_0^2} + \frac{n}{\sigma^2} \right)^{-1}$ , and  $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$ . This means the posterior mean of location parameter  $\mu$  is estimated by weighted the sample mean  $\bar{x}$  and the prior mean  $\mu_0$  via their precisions  $\sigma^2$  and  $\sigma_0^2$ . A nice interpretation of the posterior mean is that it combines information of data (sample mean) and knowledge (prior) together into the model Equation (15.5). Further, a new prediction of  $x$  given this model is also a normal distribution that

$$\hat{x} \sim N(\mu_n, \sigma_n^2 + \sigma^2). \quad (15.6)$$

In this example, the prior and the posterior are both normal distributions that we call this kind of prior as a conjugate prior. In general, a conjugate prior may not exist and may not have a good interpretation to the application. The advantage is that the analytical solution is feasible for conjugate cases. However, a prior may be better to borrow from known information such as previous experiments or domain knowledge. For instance, empirical Bayes relies on empirical data information, or non-informative priors provide wider range of parameters. Nevertheless, Markov Chain Monte Carlo (MCMC) is a typical solution when an analytical solution is tedious.

## 15.2 Hastings-Metropolis Algorithm

In reality, a proposed distribution may not be easy to obtain samples or to generate from, while Acceptant-Rejection Sampling algorithm is a fundamental method in Computational Statistics to deal with this situation by generating data from a relative easier distribution and based on the acceptant-rejection probability to keep or drop the samples. See [Ross \(1996\)](#) for more details about Acceptant-Rejection Sampling algorithm.

Hastings-Metropolis algorithm ([Hastings, 1970](#); [Metropolis et al., 1953](#)) is one of Markov Chain Monte Carlo method to obtain a sequence of random samples where a proposed distribution is difficult to sample from. The idea is to utilize Acceptant-Rejection Sampling algorithm to sample sequentially from conditional distributions provided relative easier than the proposed distribution, and via acceptance rejection probability to screen appropriate data from an equilibrium distribution. The computation of  $\pi$  (the ratio of a circle's circumference to its diameter, not prior) in Section 4.1 is an example of Acceptant-Rejection Sampling algorithm for Monte Carlo case but without Markov Chain.

Suppose a stationary distribution exists for  $\theta$  in the domain of investigation  $\Theta$ . Provided the Markov Chain is adequate (periodic, irreducible, time reversible, ...), we may have

$$\pi(\theta^{(i)})p(\theta|\theta^{(i)}) = \pi(\theta)p(\theta^{(i)}|\theta) \quad (15.7)$$

where  $p(\theta|\theta^{(i)})$  is a transition probability at the  $i$ -th step from the current state  $\theta^{(i)}$  to a new state  $\theta$  for all  $\theta^{(i)}, \theta \in \Theta$ . Since  $p(\theta|\theta^{(i)})$  may not be easy to sample, Hastings-Metropolis algorithm

suggests a proposal distribution  $q(\theta|\theta^{(i)})$  with an acceptant probability  $a(\theta|\theta^{(i)})$  such that

$$a(\theta|\theta^{(i)}) = \frac{p(\theta|\theta^{(i)})}{q(\theta|\theta^{(i)})}. \quad (15.8)$$

Equation (15.7) becomes

$$\frac{a(\theta|\theta^{(i)})}{a(\theta^{(i)}|\theta)} = \frac{\pi(\theta)q(\theta^{(i)}|\theta)}{\pi(\theta^{(i)})q(\theta|\theta^{(i)})}. \quad (15.9)$$

The acceptant probability will be

$$a(\theta|\theta^{(i)}) = \min \left\{ 1, \frac{\pi(\theta)q(\theta^{(i)}|\theta)}{\pi(\theta^{(i)})q(\theta|\theta^{(i)})} \right\} \quad (15.10)$$

that  $\theta^{(i+1)} = \theta$  if accepted, otherwise  $\theta^{(i+1)} = \theta^{(i)}$  (new  $\theta$  is rejected).

The steps of Hastings-Metropolis Algorithm are summarized next:

**Step 1:** Initial a  $\theta^{(0)}$  from  $\pi(\theta)$ . Set  $i = 1$ .

**Step 2:** Generate a new  $\theta'$  from  $g(\theta|\theta^{(0)})$ .

**Step 3:** Compute  $a(\theta'|\theta^{(i)})$ .

**Step 4:** Genera a uniform random variable  $U$ .

**Step 5:** If  $U \leq a(\theta'|\theta^{(i)})$ , then set  $\theta^{(i+1)} = \theta'$ . Otherwise, set  $\theta^{(i+1)} = \theta^{(i)}$ .

**Step 6:** Set  $i = i + 1$  and repeat Steps 2 to 5.

Typically, we repeat Steps 2 to 5 until the process is burn-in, says  $I_b = 1,000$  iterations, after that we continuously collect  $\{\theta^{(i)}\}$  for thinning every  $I_t = 10$  iterations to release time dependent problems. Repeat the thinning process until  $I_n$  samples are reached. We also repeat  $I_c = 5$  Markov Chains with different initial values to verify the stationary. The determinations of  $I_b$ ,  $I_t$ ,  $I_n$ , and  $I_c$  are dependent on models, data, and prior, see Spiegelhalter et al. (2003) for more information.

Although Hastings-Metropolis algorithm may solve complex problem, larger number of  $I_b$ ,  $I_t$ ,  $I_n$ , and  $I_c$  also result in time consuming computations and large storage space. An easy way to rescue this burden is to parallelize the algorithm. At least three possible parallelizations for  $N$  processors can be considered in following.

1. Each Markov Chain is executed on each processor. Only  $I_n/N$  samples are needed to be collected for each processor provided every Markov Chain is burn-in.
2. Execute one Markov Chain on one processor. Until the Markov Chain is burn-in, then the burn-in state is broad casted to all processors. Set different random seeds on all processors, then all processors proceed the Markov Chain until  $I_n/N$  samples are collected for each processor. Note that the approach is probably useful for short burn-in chains.
3. For large size problem, distributing data is unavoidable, then  $N$  processors execute one common Markov Chain to collect  $I_n$  samples.

We next use a galaxy velocity example to demonstrate the first parallelization above, and make statistical inference based on the Bayesian framework.

### 15.3 Galaxy Velocity

Velocities of 82 galaxies in the region of Corona Borealis are measured and reported in (Roeder, 1990), and the `galaxies` dataset is available in **MASS** package of R. The mean is about 20,828.17 km/sec and the standard deviation is about 4,563.758 km/sec. Figure 15.1 shows the distribution of data.

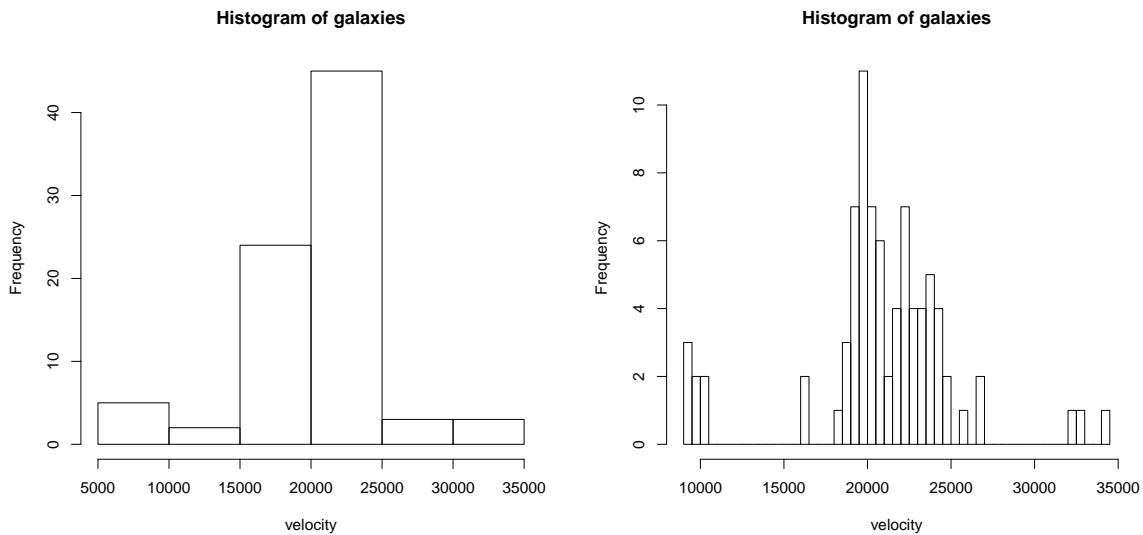


Figure 15.1: The left plot is based on default setting of `hist(galaxies)` and the right plot is based on `hist(galaxies, nclass=50)` providing more details of distribution.

Suppose we are interesting in the mean velocity of those galaxies and want to model them as Equations (15.3), (15.4), and (15.5). An example code is given in the **pbdDEMO** demo via

```
At the shell prompt, run the demo with 4 processors by
(Use Rscript.exe for windows system)
mpirun -np 4 Rscript -e "demo(mcmc_galaxy, 'pbdDEMO', ask=F, echo=F)"
```

The example has outputs next that it updates 11000 iterations in total, collects samples in every 10 iterations after 1000 burnin iterations, and 1000 samples totally collected for inference. The posterior mean of  $\mu$  (velocity of those galaxies) is about 20,820.92 km/sec and 95% credible interval is (19816.91, 21851.07) km/sec.

#### R Output

```
Total iterations: 11000
Burnin: 1000
Thinning: 10
```

```

Total samples: 1000
Posterior mean: 20820.92
95% credible interval:
 2.5% 97.5%
19816.91 21851.07

```

Also, Figure 15.2 provides the MCMC trace and posterior distribution of  $\mu$  based on those 1000 samples.

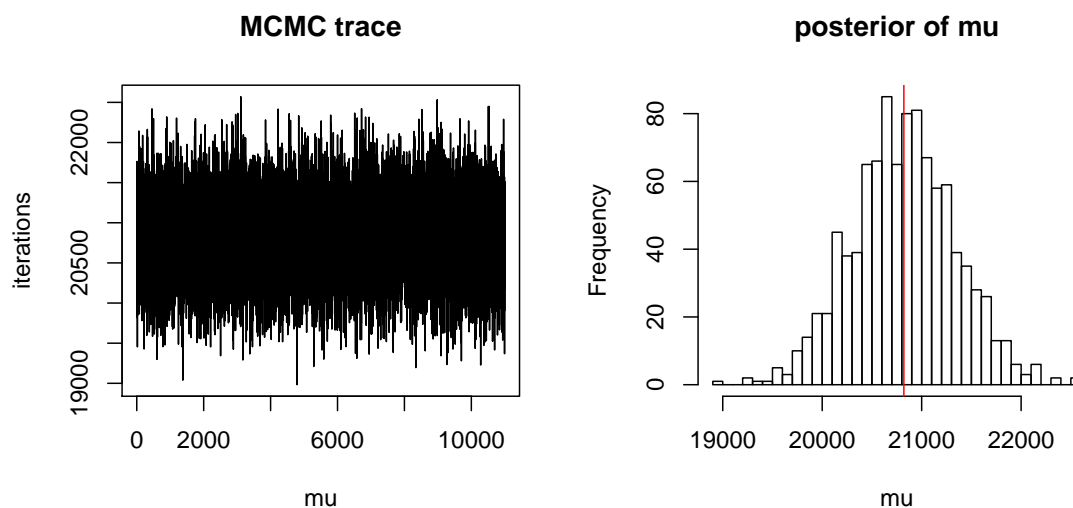


Figure 15.2: The left plot shows the MCMC trace of 11,000 iterations. The right plot displays the posterior distribution of  $\mu$  in 1000 samples and red line indicates the posterior mean.

## 15.4 Parallel Random Number Generator

We demonstrate a simple MCMC example to fit a Gaussian model to Galaxy dataset in Section 15.3, but intend to raise an important issue “parallel random number generator” to simulation technique especially in distributed environment. Even there is no need to use parallel random number generator, but the purpose here is to explain when to synchronize the random number if models are getting more complex.

**pbdMPI** builds in parallel random number generator via **rlecuyer** (Sevcikova and Rossini, 2012), and `comm.set.seed(..., diff = TRUE)` is setting different streams of random numbers to all processors. Suppose different streams provide independent (or closely) random variables, then every processor can perform MCMC on local data independently. However, new parameters and the decision for the new parameters, should be consistent in all processors so that synchronization is necessary in some stages of each MCMC iteration.

The Galaxy demo code uses 4 processors to hold the dataset and every processor start from the different seed to generate new  $\mu$  and uniform random variable  $U$  for rejection ratio, but only the

values on rank 0 are `bcast()` and used by all other ranks.

#### Hastings-Metropolis MCMC

```

1 ret <- NULL
2 ret.all <- NULL
3 mu.org <- rnorm(1, mean = mu.0, sd = sigma.0)
4 # No need to synchronize if diff = FALSE in comm.set.seed().
5 mu.org <- bcast(mu.org)
6 for(i in 1:(I.b + I.t * I.n)){
7 mu.new <- rnorm(1, mean = mu.0, sd = sigma.0)
8 # No need to synchronize if diff = FALSE in comm.set.seed().
9 mu.new <- bcast(mu.new)
10
11 a <- acceptance(x.gbd, mu.new, mu.org)
12 U <- runif(1)
13 # No need to synchronize if diff = FALSE in comm.set.seed().
14 U <- bcast(U)
15
16 if(U <= a){
17 mu.org <- mu.new
18 }
19
20 ret.all <- c(ret.all, mu.org)
21 if(i > I.b && (i %% I.t == 0)){
22 ret <- c(ret, mu.org)
23 }
24 }

```

Although we can use `comm.set.seed(..., diff = FALSE)` as default, the parameter  $\mu$  and  $U$  are tiny and in common of all ranks so that the cost of communication is relative small. For more complex models, we may consider to distribute parameters as well and make decision locally, then we can reduce the cost further. In such cases, parallel random number generators are better solutions.

## 15.5 Exercises

- 15-1 Prove Equation (15.5) and claim it is conjugate. [Hint: Equation \(15.2\).](#)
- 15-2 Prove Equation (15.6) and explain intuitively why the variance of predictive sample is increased comparing with that of observed samples. [Hint: is a 95% predictive interval wider than a 95% confidence interval.](#)
- 15-3 Claim that Equation (15.10) is the solution of Equation (15.9). [Hint: when is  \$a\(\theta^{\(i\)}|\theta\) = 1\$ ?](#)
- 15-4 Prove the proposal distribution  $q$  with Equation (15.10) provides the desired distribution



$p$ . [Hint: Acceptance-Rejection Sampling algorithm.](#)

- 15-5 Claim that the upper bound of Equation (15.8) controls the performance of Hastings-Metropolis algorithm. [Hint: what if  \$q\(\theta|\theta^{\(i\)}\) \equiv p\(\theta|\theta^{\(i\)}\)\$ ?](#)
- 15-6 Discuss when Hastings-Metropolis algorithm fails. Provide an example that is an inefficient case of Hastings-Metropolis algorithm. [Hint: What are requirements of Markov Chain?](#)
- 15-7 Extend the model and algorithm of galaxy velocities example for unknown mean and unknown variance. e.g.

$$\begin{aligned} \mathbf{x} &\stackrel{i.i.d.}{\sim} N(\mu, \sigma^2) \\ \mu &\sim N(\mu_0, \sigma_0^2) \\ \sigma &\sim \text{Gamma}(\alpha_0, \beta_0) \end{aligned}$$

Find the 95% creditable region for  $(\mu|\mathbf{x}, \sigma|\mathbf{x})$ .

- 15-8 Section 15.3 only considers homogeneous distribution for all galaxy velocities. As model-based clustering in Section 13, please extend to a two clusters problem and implement it in Bayesian framework.
- 15-9 At the end of Section 15.4, we mention a potential case to avoid communication for generating new parameters of complex MCMC models. Given an example and implement in two different ways, one uses parallel random numbers and the other uses traditional random numbers plus synchronization.
- 15-10 Formulate and implement a finite mixture model with Bayesian MCMC approach to analysis the Galaxy data. Compare results of the new implementation to the model-based clustering with EM algorithms where the Galaxy data is simply one dimensional variables in the view of Chapter 13. [Hint: Marin and Robert \(2007\).](#)
- 15-11 The trace plot of Figure 15.2 appears that the chain is converged but is considerable too “hot” that may over mixing. Try to adjust the simulation code such that acceptance rate is around 20% to 30%.

## Pairwise Distance and Comparisons

*An approximate answer to the right problem is worth a good deal more than an exact answer to an approximate problem.*

—John Tukey

### 16.1 Introduction

Distance method is not only a fundamental tool in geometry, but also appears in statistics and other applied disciplines. For example, least square method in regression can be simply derived and computed via Euclidean distance. The resulting line is an approximate answer in terms of minimum total distance to all observations. Distance is also related to a similarity measure of two observations describing relationship of the two. Usually, the smaller of distance the closer of relation. For example, the higher probability (probability is a measure) of one virus evolving to a mutant means the smaller distance (related closely) of two viruses as described in Chapter 14. Further, distance method is simple to apply on clustering problems and easy to visualize data structures such as K-means algorithm which is a special case of model-based clustering introduced in Chapter 13. For instance, the observations of the same group are more similar in characteristics with each other than those between different groups.

Potentially, computing distance of several observations involves half of pairwise comparisons if distance is symmetric, and involves all pairwise comparisons if distance is not symmetric. Also, if number of observations is small, then most of distance methods can be compute efficient within one core. For moderate number of observations or complex distance systems, the computing can be parallelized wisely in several levels. For example, one may utilize multiple threads or co-processors to archive performance gains. For large number of observations, the computing is not trivial if data are distributed across cores. Further, the dimension of resulting distance array may be much larger the number of observations and may only be held distributed across cores. Note that for some models or iterative algorithms, it is not wise to dump the distance array into disk since that decreases performance due to overhead cost for I/O. For example, one may utilize distributed parallelization to avoid these restrictions.

In the context of `pbdR`, we focus on distributed methods and abstract computing of distance to allow user-defined comparison (dissimilarity) functions of any two observations. We briefly introduce issues and methods of distributed distance and comparisons first, and followed by demonstration of hierarchical clusterings on the `iris` dataset of Chapter 13. This example can be done using exists distance function in R. Further, we provide a biological application of building phylogenetic trees on the *Pony 524* dataset of Chapter 14 utilizing evolutionary models to compute probability distance. This example demonstrate how user-defined function can be defined and used to obtain special distance. In general, the function can be extended to multiple comparisons and tests.

## 16.2 Distributed Distance and Comparisons

Suppose  $x$  and  $y$  are two observations and  $d(x, y)$  is a distance or a comparison of  $x$  and  $y$ . Note that  $x$ ,  $y$ , and  $d(\cdot, \cdot)$  could be very generic as long as they are well defined. Although, it is efficient to compute a distance of any two observations in R via `dist()` serially, it becomes non-trivial to compute distance of distributed observations in parallel.

The potential problems include:

- (P1) Communication must be evoked between processors when any two observations are not located within the same processor.
- (P2) The resulting distance matrix may be too big to held in one processor as data size increased even only a half (lower triangular matrix is stored as row-major in a vector.)
- (P3) Compute all comparisons may be too time consuming even for small data sets.

Distributed situations of observations and computed results (distance matrix) are categorized next.

- (C1) Both observations and distance matrix are in one node and may both be in serial or in parallel within the node, typically via OpenMP ([OpenMP ARB, 1997](#)).
- (C2) Observations are in common in all processors and distance matrix is distributed across nodes.
- (C3) Observations are distributed across nodes and distance matrix is in common in all nodes.
- (C4) Both observations and distance matrix are distributed across nodes.

Here, we may presume the distribution method is GBD row-major matrix (or row-block major) as introduced in Section 3.3 since most of native R functions can be extended and reused in such a similar way.

Note that the `dist()` only supports a few distance methods and assume distance is symmetric by definition. However, in practice, a more general measure may not be necessarily symmetric of two observations. i.e.  $d(x, y) \neq d(y, x)$ . In some cases,  $d(x, x) \neq 0$  and the distance may also be dependent on other measurements or conditions. In general, a function for comparing any two  $x$  and  $y$  is possible to replace `dist()`.

### 16.3 Hierarchical Clustering

Hierarchical clustering is a popular statistical tool in fundamental multivariate statistics and is heavily relied on a distance matrix to classify data. Several algorithms are proposed to build dendrograms or trees, then prune branches of the resulting trees to identify possible subgroups. The basic function `hclust()` takes a dissimilarity structure as produced by `dist()` and returns a tree object can be visualized. The method option “average” linkage is equivalent to UPGMA (Unweighted Pair Group Method with Arithmetic Mean) method (Sokal and Michener, 1958) which is one of popular methods in ecology for classification.

For example, the `iris` dataset used in Chapter 13 can be clustered in hierarchical clustering. First, we distribute 150 observations in four cores and compute Euclidean distances in four dimensional space (‘Sepal.Length’, ‘Sepal.Width’, ‘Petal.Length’, and ‘Petal.Width’). Note that the distance may not be meaningful to the data, but preserve some (dis-) similarity of the observations. We compute the dissimilarity matrix in distributed manners via a utility function `comm.dist()` of `pbdMPI` (Chen et al., 2012a) and store the result in a common matrix across all cores. We based on the matrix to perform a UPGMA clustering. The example in SPMD can be found in demo via

```
At the shell prompt, run the demo with 4 processors by
(Use Rscript.exe for windows system)
mpirexec -np 4 Rscript -e "demo(dist_iris,'pbdDEMO',ask=F,echo=F)"
```

and it returns a dendrogram as Figure 16.1 where species “Versicolor” (in green) and “Virginica” (in blue) are potentially overlapped and differ from “Setosa” (in red).

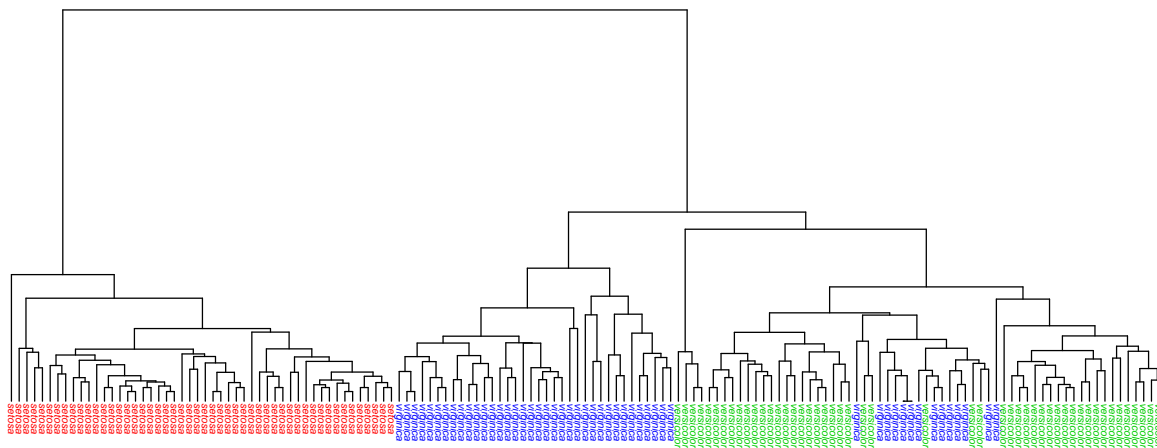


Figure 16.1: Hierarchical clustering result of `iris` dataset.

## 16.4 Neighbor Joining

In some sense, Figure 16.1 is a rooted tree and the “average” method as well as UPGMA assumes a constant rate of evolution (molecular clock hypothesis). However, these assumption may not be appropriate to most sequence evolutionary topics where a gene tree should be more suitable to interpret relation of sequences or species. We introduce a popular approach in evolution biology and build a evolutionary tree for *Pony 524* dataset. We select JC69 evolutionary model (Jukes and Cantor, 1969) as a probability measure to compute for distance (evolution time) of 146 EIAV sequences and use a neighbor joining tree (Saitou and Nei, 1987) to build an unrooted tree.

The purpose is to design a wrapper function, says `my.dist(x, y)`, that takes a pairs of sequences `x` and `y` as inputs, and returns a user-defined distance of given data. The utility function `comm.pairwise()` of `pbdMPI` (Chen et al., 2012a) is more flexible than `comm.dist()`. Through the options `pairid.gbd` and `FUN = my.dist`, the function can evaluate `my.dist()` on the given dataset `X` in row major blocks. For *Pony 524*, the `X` is the DNA sequences and `my.dist()` is a wrapper of `phyclust.edist`.

The example in SPMD can be found in demo via

```
At the shell prompt, run the demo with 4 processors by
(Use Rscript.exe for windows system)
mpiexec -np 4 Rscript -e "demo(dist_pony,'pbdDEMO',ask=F,echo=F)"
```

and it returns a neighbor-joining tree as Figure 16.2.

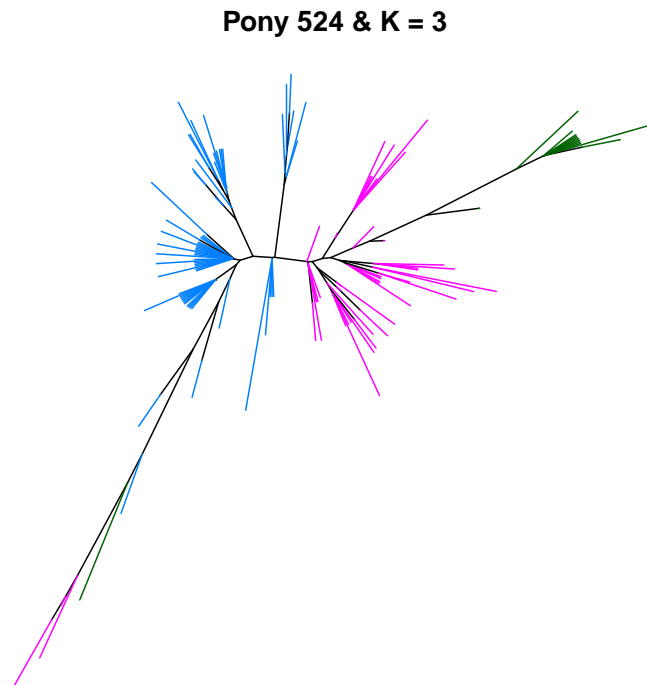


Figure 16.2: Neighbor-joining tree of *Pony 524* dataset colored by three clusters.

## 16.5 Exercises

- 16-1 What are potential limitations of distance approaches?
- 16-2 Prove that clustering based on Euclidean distance is equivalent to that clustering based on multivariate Normal distributions with identity variance covariance matrices.
- 16-3 Prove that the “average” method of `hclust()` is equivalent to the UPGMA method.
- 16-4 Given  $n$  observations or taxa, analytically find total numbers of possible rooted and unrooted trees,  $(2n - 5)!!$  and  $(2n - 3)!!$ , respectively.
- 16-5 As number of observations increases, the data and the distance matrix are both distributed as the category (C4). State potential problems of implementations and minimum costs of communications.
- 16-6 Discuss the difficulties and problems of designing tree algorithms on a distributed manner.

Part VI

Appendix



## Numerical Linear Algebra and Linear Least Squares Problems

*Mathematics is written for mathematicians.*

—Nicolaus Copernicus

For the remainder, assume that all matrices are real-valued.

Let us revisit the problem of solving linear least squares problems, introduced in Section 4.5. Recall that we wish to find a solution  $\beta$  such that

$$\|X\beta - y\|_2^2$$

In the case that  $X$  is full rank (which is often assumed, whether reasonable or not), this has analytical solution

$$\hat{\beta} = (X^T X)^{-1} X^T y \quad (\text{A.1})$$

However, even with this nice closed form, implementing this efficiently on a computer is not entirely straightforward. Herein we discuss several of the issues in implementing the linear least squares solution efficiently. For simplicity, we will assume that  $X$  is full rank, although this is not necessary — although rank degeneracy does complicate things. For more details on the rank degeneracy problem, and linear least squares problems in general, see the classic *Matrix Computations* (Golub and Van Loan, 1996).

### A.1 Forming the Normal Equations

If we wish to implement this numerically, then rather than directly computing the inverse of  $X^T X$  directly, we would instead compute the Cholesky factorization

$$X^T X = LL^T$$

where  $L$  is lower triangular. Then turning to the so-called “normal equations”

$$(X^T X) \beta = X^T y \quad (\text{A.2})$$



by simple substitution and grouping, we have

$$L(L^T\beta) = \mathbf{X}^T\mathbf{y}$$

Now, since  $L$  is triangular, these two triangular systems (one forward and one backward substitution found by careful grouping of terms above) can be solved in a numerically stable way (Higham, 2002). However, forming the Cholesky factorization itself suffers from the effects of roundoff error in having to form the product  $\mathbf{X}^T\mathbf{X}$ . We elaborate on this to a degree in the following section.

## A.2 Using the QR Factorization

Directly computing the normal equations is ill advised, because it is often impossible to do so with adequate numerical precision. To fully appreciate this problem, we must entertain a brief discussion about condition numbers.

By definition, if a matrix has finite condition number, then it must have been invertible. However, for numerical methods, a condition number which is “big enough” is essentially infinite (loosely speaking). And observe that forming the product  $\mathbf{X}^T\mathbf{X}$  squares the condition number of  $\mathbf{X}$ :

$$\begin{aligned}\kappa(\mathbf{X}^T\mathbf{X}) &= \|\mathbf{X}^T\mathbf{X}\| \left\| (\mathbf{X}^T\mathbf{X})^{-1} \right\| \\ &= \|\mathbf{X}^T\mathbf{X}\| \left\| \mathbf{X}^{-1} (\mathbf{X}^T)^{-1} \right\| \\ &= \|\mathbf{X}^T\| \|\mathbf{X}\| \|\mathbf{X}^{-1}\| \|\mathbf{X}^{-T}\| \\ &= \|\mathbf{X}\| \|\mathbf{X}\| \|\mathbf{X}^{-1}\| \|\mathbf{X}^{-1}\| \\ &= \|\mathbf{X}\|^2 \|\mathbf{X}^{-1}\|^2 \\ &= \kappa(\mathbf{X})^2\end{aligned}$$

So if  $\kappa(\mathbf{X})$  is “large”, then forming this product can lead to large numerical errors when attempting to numerically invert or factor a matrix, or solve a system of equations.

To avoid this problem, the orthogonal QR Decomposition is typically used. Here we take

$$\mathbf{X} = \mathbf{Q}\mathbf{R}$$

where  $\mathbf{Q}$  is orthogonal and  $\mathbf{R}$  is upper trapezoidal (in the overdetermined case,  $\mathbf{R}$  is triangular). This is beneficial, because orthogonal matrices are norm-preserving, i.e.  $\mathbf{Q}$  is an isometry, and whence

$$\begin{aligned}\|\mathbf{X}\beta - \mathbf{y}\|_2 &= \|\mathbf{Q}\mathbf{R}\beta - \mathbf{y}\|_2 \\ &= \|\mathbf{Q}^T\mathbf{Q}\mathbf{R}\beta - \mathbf{Q}^T\mathbf{y}\|_2 \\ &= \|\mathbf{R}\beta - \mathbf{Q}^T\mathbf{y}\|_2\end{aligned}$$

This amounts to solving the triangular system

$$R\beta = Q^T \mathbf{y}$$

As noted in Section A.1, solving this system can be done in a numerically stable fashion (and the high performance libraries employed by both R and **python** use stable implementations). The key difference here is that the QR factorization is of  $\mathbf{X}$ , not  $\mathbf{X}^T \mathbf{X}$ , and so we need only worry about the conditioning of  $\mathbf{X}$  (as opposed to its squared condition number).

While this method is much less prone to the numerical issues discussed above, it is much slower computationally. Additionally, we note that unlike the method in forming the normal equations, this method can be extended to the rank degenerate case.

### A.3 Using the Singular Value Decomposition

There is another, arguably much more well-known matrix factorization which we can use to develop yet another analytically equivalent solution to the least squares problem, namely the singular value decomposition (SVD). Using this factorization leads to a very elegant solution, as is so often the case with the SVD.

Note that in (A.1), the quantity

$$(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$$

is the Moore-Penrose inverse of  $\mathbf{X}$ . So if we take

$$\mathbf{X} = U\Sigma V^T$$

to be the SVD of  $\mathbf{X}$ , then we have

$$\begin{aligned} \mathbf{X}^+ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \\ &= \left( (U\Sigma V^T)^T (U\Sigma V^T) \right)^{-1} U\Sigma V^T \\ &= (V\Sigma^T \Sigma V^T)^{-1} V\Sigma^T U^T \\ &= V \left( (\Sigma^T \Sigma)^{-1} \Sigma^T \right) U^T \\ &= V\Sigma^+ U^T \end{aligned}$$

Whence,

$$\beta = V\Sigma^+ U^T \mathbf{y}$$

Conceptually, this is arguably the most elegant method of solving the linear least squares problem. Additionally, as with the QR method above, with slight modification the above argument can extend to the rank degenerate case; however, we suspect that the SVD is much more well known to mathematicians and statisticians than is the QR decomposition. This abstraction comes at a great cost, though, as this approach is handily the most computationally intensive of the three presented here.



## Linear Regression and Rank Degeneracy in R

*When a doctor does go wrong, he is the first of criminals. He has the nerve and he has the knowledge.*

—Sherlock Holmes

In the case that  $\mathbf{X}$  is rank deficient, then  $\mathbf{X}$  (and whence  $\mathbf{X}^T \mathbf{X}$ ) is not invertible, so the problem can not be solved by the method of Section A.1. Both R and **pbdr** use a QR factorization as in Section A.2, although the two systems use a slightly different approach. While most of the linear algebra in R is handled by LAPACK (Anderson et al., 1999), arguably the most important numerical function in all of R, namely `lm.fit()` used by `lm()` to fit linear regression models, uses LINPACK (Dongarra et al., 1979). By comparison to LAPACK, LINPACK is much less sophisticated. However, **pbdr** uses level 3 PBLAS and ScaLAPACK (the distributed equivalent of using level 3 BLAS and LAPACK) to fit linear regression models.

The LINPACK routines used by R are DQRSL, which calls a modified DQRDC2 to compute a rank-revealing QR factorization with a “limited pivoting strategy” (more on this later). Finally, DQRSL is called to apply the output of the QR factorization to compute the least squares solutions. By contrast, **pbdr** uses a modified PDGELS routine, which uses a version of PDGEQPF modified to use R’s “limited pivoting strategy”, and then calls PDORMQR to fit the least squares solution.

Neither approach assumes that the model matrix is full rank. Instead, the methods are *rank-revealing*, in that they attempt to discover the numerical rank while computing the orthogonal factorization. However, both R and (for the sake of consistency) **pbdr** use a “limited pivoting strategy” (with language, we believe, due to Ross Ihaka) in determining numerical rank. Generally when computing a QR with pivoting, for the sake of numerical stability one chooses the column with largest partial norm while forming the Householder reflections. However, in doing so it is possible to permute the columns in such a way that a desired statistical interpretation (such as in an ANOVA) is destroyed. The solution employed by R is to merely iterate over the columns and choose the current column as the pivot each time. When a column is detected to have “small” partial norm, it is pushed to the back. The author of these modification writes:

a limited column pivoting strategy based on the 2-norms of the reduced columns moves columns with near-zero norm to the right-hand edge of the x matrix. this strategy means that sequential one degree-of-freedom effects can be computed in a natural way.

i am very nervous about modifying linpack code in this way. if you are a computational linear algebra guru and you really understand how to solve this problem please feel free to suggest improvements to this code.

So in this way, if a model matrix is full rank, then the estimates coming from R should be considered at least as trustworthy as probably every other statistical software package of note. If it is not, then this method presents a possible numerical stability issue; although to what degree, if any at all, this is actually a problem, the authors at present have no real knowledge. If numerical precision is absolutely paramount, consider using the SVD to solve the least squares problem; though do be aware that this is hands down the slowest possible approach.

We again note that the limited pivoting strategy of R is employed by `pbdR` in the `lm.fit()` method for class `ddmatrix`.

## Part VII

# Miscellany

## References

- D. Adler, C. Gläser, O. Nenadic, J. Oehlschlägel, and W. Zucchini. ff: memory-efficient storage of large data on disk and fast access functions, 2013. URL <http://CRAN.R-project.org/package=ff>. R Package version 2.2-11.
- H. Akaike. A new look at the statistical model identification. *IEEE Transaction on Automatic Control*, 19:716–723, 1974.
- Revolution Analytics. *foreach: Foreach looping construct for R*, 2012. URL <http://CRAN.R-project.org/package=foreach>. R Package version 1.4.0.
- E. Anderson, Z. Bai, C. Bischof, L. S. Blackford, J. Demmel, Jack J. Dongarra, J. Du Croz, S. Hammarling, A. Greenbaum, A. McKenney, and D. Sorensen. *LAPACK Users' guide (third ed.)*. Society for Industrial and Applied Mathematics, Philadelphia, PA, USA, 1999. ISBN 0-89871-447-8.
- P. Baccam, R.J. Thompson, Y. Li, W.O. Sparks, M. Belshan, K.S. Dorman, Y. Wannemuehler, J.L. Oaks, J.L. Cornette, and S. Carpenter. Subpopulations of equine infectious anemia virus rev coexist in vivo and differ in phenotype. *Journal of Virology*, 77(22):12122–12131, 2003.
- F. Benson. A note on the estimation of mean and standard deviation from quantiles. *Journal of the Royal Statistical Society. Series B (Methodological)*, 11(1):91–100, 1949.
- J.O. Berger. *Statistical Decision theory and Bayesian Analysis*. Springer, 2nd edition, 1993.
- L. S. Blackford, J. Choi, A. Cleary, E. D'Azevedo, J. Demmel, I. Dhillon, J. Dongarra, S. Hammarling, G. Henry, A. Petitet, K. Stanley, D. Walker, and R. C. Whaley. *ScaLAPACK Users' Guide*. Society for Industrial and Applied Mathematics, Philadelphia, PA, 1997. ISBN 0-89871-397-8 (paperback). URL [http://netlib.org/scalapack/slug/scalapack\\_slug.html/](http://netlib.org/scalapack/slug/scalapack_slug.html/).
- S. Carpenter, W.-C. Chen, and K.S. Dorman. Rev variation during persistent lentivirus infection. *Viruses*, 3:1–11, 2011.
- G. Casella and R.L. Berger. *Statistical Inference*. Cengage Learning, 2nd edition, 2001.

- W.-C. Chen. Overlapping codon model, phylogenetic clustering, and alternative partial expectation conditional maximization algorithm. *Ph.D. Diss., Iowa Stat University*, 2011.
- W.-C. Chen and R. Maitra. Model-based clustering of regression time series data via APECM — an AECM algorithm sung to an even faster beat. *Statistical Analysis and Data Mining*, 4: 567–578, 2011.
- W.-C. Chen and G. Ostrouchov. Hpsc – high performance statistical computing for data intensive research, 2011. URL <http://thirteen-01.stat.iastate.edu/snoweye/hpsc/>.
- W.-C. Chen and G. Ostrouchov. pmclust: Parallel model-based clustering, 2012. URL <http://cran.r-project.org/package=pmclust>. R Package.
- W.-C. Chen, G. Ostrouchov, D. Schmidt, P. Patel, and H. Yu. pbdMPI: Programming with big data – interface to MPI, 2012a. URL <http://cran.r-project.org/package=pbdMPI>. R Package.
- W.-C. Chen, G. Ostrouchov, D. Schmidt, P. Patel, and H. Yu. *A Quick Guide for the pbdMPI package*, 2012b. URL <http://cran.r-project.org/package=pbdMPI>. R Vignette.
- W.-C. Chen, D. Schmidt, G. Ostrouchov, and P. Patel. pbdSLAP: Programming with big data – scalable linear algebra packages, 2012c. URL <http://cran.r-project.org/package=pbdSLAP>. R Package.
- W.-C. Chen, G. Ostrouchov, D. Pugmire, M. Prabhat, and M. Wehner. A parallel em algorithm for model-based clustering with application to explore large spatio-temporal data. *Technometrics*, 55:513–523, 2013.
- Wei-Chen Chen, Ranjan Maitra, and Volodymyr Melnykov. EMCluster: EM algorithm for model-based clustering of finite mixture gaussian distribution, 2012d. R Package, URL <http://cran.r-project.org/package=EMCluster>.
- A.P. Dempster, N.M. Laird, and D.B. Rubin. Maximum likelihood for incomplete data via the EM algorithm (with discussion). *Journal of the Royal Statistical Society, Series B*, 39:1–38, 1977.
- J. J. Dongarra, C. B. Moler, J. R. Bunch, and G. W. Stewart. *LINPACK User’s Guide*. SIAM, 1979.
- B. Eaton. *User’s Guide to the Community Atmosphere Model CAM-5.1*. NCAR, 2011. URL <http://www.cesm.ucar.edu/models/cesm1.0/cam/>.
- B. Efron. Bootstrap methods: Another look at the jackknife. *Annals of Statistics*, 1979.
- Thomas S. Ferguson. *A Course in Large Sample Theory*. Chapman & Hall/CRC, 1996. ISBN 978-0412043710.
- R. A. Fisher. On the Mathematical Foundations of Theoretical Statistics. *Philosophical Transactions of the Royal Society of London. Series A, Containing Papers of a Mathematical or Physical Character*, 222:309–368, 1922.
- R.A. Fisher. The use of multiple measurements in taxonomic problems. *Annals of Eugenics*,

2:179–188, 1936.

E. Forgy. Cluster analysis of multivariate data: efficiency vs. interpretability of classifications. *Biometrics*, 21:768–780, 1965.

C. Fraley and A.E. Raftery. Model-based clustering, discriminant analysis, and density estimation. *Journal of the American Statistical Association*, 97:611–631, 2002.

C. Fraley, A. Raftery, and L. Scrucca. mclust: Normal mixture modeling for model-based clustering, classification, and density estimation, 1999. R Package, URL <http://cran.r-project.org/package=mclust>.

Benjamin Franklin. Join, or die [cartoon]. In *The Library of Congress*. Retrieved from <http://www.loc.gov/pictures/item/2002695523>, 1754.

A. Gelman, J.B. Carlin, H.S. Stern, and D.B. Rubin. *Bayesian Data Analysis*. Chapman & Hall/CRC, 2nd edition, 2003.

Gene H. Golub and Charles F. Van Loan. *Matrix Computations (Johns Hopkins Studies in Mathematical Sciences)(3rd Edition)*. The Johns Hopkins University Press, 3rd edition, October 1996.

W. Gropp, E. Lusk, and A. Skjellum. *Using MPI: Portable Parallel Programming with the Message-Passing Interface*. Cambridge, MA, USA: MIT Press Scientific And Engineering Computation Series, 1994.

W.K. Hastings. Monte carlo sampling methods using markov chains and their applications. *Biometrika*, 57:97–109, 1970.

Nicholas J. Higham. *Accuracy and Stability of Numerical Algorithms*. Society for Industrial and Applied Mathematics, Philadelphia, PA, USA, 2nd edition, 2002. ISBN 0898715210.

L. Hubert and P. Arabie. Comparing partitions. *Journal of Classification*, 2:193–218, 1985.

R.R. Hudson. Generating samples under a wright-fisher neutral model of genetic variation. *Bioinformatics*, 18:337–338, 2002.

J.L.W.V. Jensen. Sur les fonctions convexes et les inégalités entre les valeurs moyennes. *Acta Mathematica*, 30:175–193, 1906.

T.H. Jukes and C.R. Cantor. *Evolution of Protein Molecules*. New York: Academic Press, 1969.

M.J. Kane and J.W. Emerson. The bigmemory project, 2010. URL <http://www.bigmemory.org>.

K. Kunen. *Set Theory: An Introduction to Independence Proofs*. North-Holland, 1980.

C. Leroux, J.-J. Cadoré, and R.C. Montelaro. Equine infectious anemia virus (eiv): What has hiv’s country cousin got to tell us? *Veterinary Research*, 35:485–512, 2004.

J.M. Marin and C. Robert. *Bayesian Core: A Practical Approach to Computational Bayesian Statistics*. Springer Texts in Statistics. Springer, 2007. ISBN 9780387389790. URL <https://>



[//www.ceremade.dauphine.fr/~xian/BCS/](http://www.ceremade.dauphine.fr/~xian/BCS/).

P. McCullagh and J.A. Nelder. *Generalized Linear Models*. Chapman & Hall/CRC, 2nd edition, 1989. ISBN 978-0412317606.

V. Melnykov and R. Maitra. Finite mixture models and model-based clustering. *Statistics Surveys*, 4:80–116, 2010.

Volodymyr Melnykov, Wei-Chen Chen, and Ranjan Maitra. MixSim: An R package for simulating data to study performance of clustering algorithms. *Journal of Statistical Software*, 51(12):1–25, 2012. URL <http://www.jstatsoft.org/v51/i12/>.

X.L. Meng and D. van Dyk. The EM algorithm — an old folk-song sung to a fast new tune (with discussion). *Journal of the Royal Statistical Society B*, 59:511–567, 1997.

N. Metropolis, A.W. Rosenbluth, M.N. Rosenbluth, A.H. Teller, and E. Teller. Equations of state calculations by fast computing machines. *Journal of Chemical Physics*, 1953.

NetCDF Group. Network common data form, 2008. URL <http://www.unidata.ucar.edu/software/netcdf/>. Software package.

OpenMP ARB. Openmp, 1997. URL <http://www.openmp.org/>.

G. Ostrouchov, W.-C. Chen, D. Schmidt, and P. Patel. Programming with big data in R, 2012. URL <http://r-pbd.org/>.

P. Patel, G. Ostrouchov, W.-C. Chen, D. Schmidt, and D. Pierce. pbdNCDF4: Programming with big data – interface to parallel unidata NetCDF4 format data files, 2013a. URL <http://cran.r-project.org/package=pbdNCDF4>. R Package.

P. Patel, G. Ostrouchov, W.-C. Chen, D. Schmidt, and D. Pierce. A quick guide for the pbd-NCDF4 package, 2013b. URL <http://cran.r-project.org/package=pbdNCDF4>. R Vignette.

David Pierce. ncdf4: Interface to unidata netcdf (version 4 or earlier) format data files, 2012. URL <http://CRAN.R-project.org/package=ncdf4>. R Package.

R Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria, 2012a. URL <http://www.r-project.org/>. ISBN 3-900051-07-0.

R Core Team. parallel: Support for parallel computation in R, 2012b. R Package.

A. Rambaut and N.C. Grassly. Seq-gen: An application for the monte carlo simulation of dna sequence evolution along phylogenetic trees. *Comput Appl Biosci*, 13(3):235–238, 1997.

K. Roeder. Density estimation with confidence sets exemplified by superclusters and voids in the galaxies. *Journal of the American Statistical Association*, 85:617–624, 1990.

S.M. Ross. *Simulation*. Oxford, 2nd edition, 1996.

N. Saitou and M. Nei. The neighbor-joining method: a new method for reconstructing phylogenetic trees. *Molecular Biology and Evolution*, 4:406–425, 1987.

- D. Schmidt, W.-C. Chen, G. Ostrouchov, and P. Patel. pbdBASE: Programming with big data – core pbd classes and methods, 2012a. URL <http://cran.r-project.org/package=pbdBASE>. R Package.
- D. Schmidt, W.-C. Chen, G. Ostrouchov, and P. Patel. *A Quick Guide for the pbdBASE package*, 2012b. URL <http://cran.r-project.org/package=pbdBASE>. R Vignette.
- D. Schmidt, W.-C. Chen, G. Ostrouchov, and P. Patel. pbdDMAT: Programming with big data – distributed matrix algebra computation, 2012c. URL <http://cran.r-project.org/package=pbdDMAT>. R Package.
- D. Schmidt, W.-C. Chen, G. Ostrouchov, and P. Patel. *A Quick Guide for the pbdDMAT package*, 2012d. URL <http://cran.r-project.org/package=pbdDMAT>. R Vignette.
- D. Schmidt, G. Ostrouchov, W.-C. Chen, and P. Patel. Tight coupling of R and distributed linear algebra for high-level programming with big data. In Patrick Kellenberger, editor, *2012 SC Companion: High Performance Computing, Networking Storage and Analysis*. IEEE Computer Society, 2012e.
- D. Schmidt, W.-C. Chen, G. Ostrouchov, and P. Patel. pbdDEMO: Programming with big data – demonstrations of pbd packages, 2013. URL <http://cran.r-project.org/package=pbdDEMO>. R Package.
- H. Sevcikova and T. Rossini. *rlecuyer: R interface to RNG with multiple streams*, 2012. URL <http://CRAN.R-project.org/package=rlecuyer>. R Package version 0.6-1.
- R. Sokal and C. Michener. A statistical method for evaluating systematic relationships. *University of Kansas Science Bulletin*, 38:1409–1438, 1958.
- D. Spiegelhalter, A. Thomas, N. Best, and D. Lunn. *WinBUGS User Manual*, 2003. URL <http://www.mrc-bsu.cam.ac.uk/bugs/winbugs/manual14.pdf>.
- Luke Tierney, A. J. Rossini, Na Li, and H. Sevcikova. snow: Simple network of workstations, 2012. URL <http://cran.r-project.org/package=snow>. R Package (v:0.3-9).
- S. Urbanek. *multicore: Parallel processing of R code on machines with multiple cores or CPUs*, 2011. URL <http://CRAN.R-project.org/package=multicore>. R Package version 0.1-7.
- M. Vertenstein, T. Craig, A. Middleton, D. Feddema, and C. Fischer. *CESM1.0.4 User's Guide*. NCAR, 2011. URL <http://www.cesm.ucar.edu/models/cesm1.0/cesm/>.
- R.A. Weiss. The discovery of endogenous retroviruses. *Retrovirology*, 3:67, 2006.
- Steve Weston. *doMPI: Foreach parallel adaptor for the Rmpi package*, 2010. URL <http://CRAN.R-project.org/package=doMPI>. R Package version 0.1-5.
- Hadley Wickham. *ggplot2: elegant graphics for data analysis*. Springer New York, 2009. ISBN 978-0-387-98140-6. URL <http://had.co.nz/ggplot2/book>.
- Hao Yu. Rmpi: Parallel statistical computing in R. *R News*, 2(2):10–14, 2002. URL [http://cran.r-project.org/doc/Rnews/Rnews\\_2002-2.pdf](http://cran.r-project.org/doc/Rnews/Rnews_2002-2.pdf).

adjusted Rand index, 97  
 AIC, 107  
 Algorithm  
   Acceptant-Rejection Sampling, 112  
   AECM, 94  
   APECM, 94  
   APECMa, 94  
   EGM, 94  
   EM, 93, 104  
   Hastings-Metropolis, 112  
   k-means, 94, 102  
   neighbor-joining, 121  
 ARI, 97  
 bootstrap, 107  
 Class  
   gbd, 26  
   .dmat, 15  
   .gbd, 15, 25, 41  
   ddmatrix, 15, 54, 57, 77, 79, 80, 84, 95, 128  
   ncvar4, 75  
   dx, 66, 69  
 Code  
   La.svd(), 63  
   RRand(), 97  
   allgather(), 25, 26  
   allreduce(), 25, 26, 36, 38  
   as.ddmatrix(), 55, 57  
   barrier(), 24  
   bcast(), 25  
   chol(), 63  
   comm.all(), 30, 31  
   comm.any(), 31  
   comm.cat(), 26  
   comm.dist(), 120  
   comm.pairwise(), 121  
   comm.print(), 26  
   comm.rank(), 24  
   comm.set.seed(), 26  
   comm.size(), 24  
   comm.sort(), 31  
   ddmatrix(), 55  
   ddmatrix, 15  
   dist(), 119  
   dmat2gbd(), 84  
   finalize(), 24  
   gather(), 25  
   gbd2dmat(), 84  
   get.jid(), 15  
   hclust(), 120  
   init(), 24  
   lm.fit(), 41  
   load.balance(), 83  
   lu(), 63  
   mclapply(), 16  
   nlm(), 40  
   optim(), 40, 92  
   pbdApply(), 28  
   pbdLapply(), 28

- `pbdsapply()`, 28
- `phyclust.edist()`, 121
- `prcomp()`, 69
- `qr()`, 63
- `quantile()`, 39
- `reduce()`, 25
- `svd()`, 63
- `uniroot()`, 40
- `unload.balance()`, 83
- conjugate prior, 112
- continuous time Markov chain, 103
- Convergence
  - a.s., 35
  - almost surely, 42
  - in distribution, 42
  - in probability, 42
- CSV, 72
- CTMC, 103
- Data
  - iris, 95, 119
  - Pony 524, 106, 119, 121
  - TREFHT, 74
- Decomposition
  - Cholesky, 63
  - eigenvalues decomposition, 109
  - LU, 63
  - QR, 41, 127
  - SVD, 12, 126
- Distribution
  - chi-squared, 90
  - multivariate normal distribution, 88, 90, 91
  - MVN, 90, 91
  - normal distribution, 39, 89
- EIAV, 105
- Gauss-Markov Theorem, 42, 89
- Gaussian mixture model, 93
- GBD column-major matrix, 21, 86
- GBD data structure, 21
- GBD row-major matrix, 21, 85
- general block distributed, 15
- HPSC, 34
- i.i.d., 88
- Library
  - BLACS, 49, 55, 59, 78–81, 84
  - Hadoop, 14
  - HDF5, 74
  - LAPACK, 127
  - LINPACK, 127
  - MPICH2, 9, 18
  - MPT, 18
  - NetCDF4, 3, 74
  - OpenMP, 16, 119
  - OpenMPI, 9, 18
  - PBLAS, 127
  - ScaLAPACK, 3, 46
- likelihood function, 88
- likelihood ratio test, 90
- linear mixed effect models, 41
- LRT, 90
- MCMC, 112
- Message Passing Interface, *see* MPI
- MLE, 89
- Model-Based Clustering, 93
- Monte Carlo, 107
- MPI, 18, 34, 74
- OLS, 40
- ordinary least squares, 40
- Package
  - EMCluster**, 94
  - MASS**, 114
  - MixSim**, 95, 97
  - Rmpi**, 12, 19, 94
  - bigmemory**, 73
  - doMPI**, 19
  - ff**, 73
  - foreach**, 19
  - mclust**, 94
  - multicore**, 16
  - ncdf4**, 74
  - parallel**, 12, 32
  - pmclust**, 94, 97
  - rlecuyer**, 26, 59, 115
  - snow**, 12
- Parallelism
  - embarrassingly parallel, 12
  - forking, 16
  - loosely coupled, 12

- manager/workers paradigm, [12](#), [14](#)
- manager/works paradigm, [108](#)
- MapReduce, [14](#)
- multi-threading, [16](#)
- SPMD, [3](#), [12](#), [14](#), [19](#), [24](#), [94](#), [107](#)
- task parallelism, [12](#), [107](#)
- tightly coupled, [12](#)
- PCA, [68](#), [92](#)
- pdf, [88](#)
- phyloclustering, [103](#)
- Principal Components Analysis, *see* PCA, *see* PCA
- re-sampling technique, [107](#)
- RNG
  - Parallel Random Number Generator, [115](#)
- semi-supervised learning, [94](#), [102](#)
- Single Program/Multiple Data, *see* SPMD
- singular value decomposition, [12](#), [126](#)
- SLLN, [35](#)
- Strong Law of Large Numbers, [35](#)
- unsupervised learning, [93](#)
- UPGMA, [120](#)
- Weak Law of Large Numbers, [42](#)
- weighted least squares, [41](#)
- WLLN, [42](#)
- WLS, [41](#)