

并行机器编程

Norm Matloff (著)

加州大学戴维斯分校

寇强 (译)

华南统计科学研究中心

印第安纳大学

新浪微博: Gossip_useR

GPU、多核、集群

本书使用 Creative Commons license 发布

<http://heather.cs.ucdavis.edu/~matloff/probstatbook.html>

本书更新网址:

英文版: <http://heather.cs.ucdavis.edu/~matloff/158/PLN/ParProcBook.pdf>

中文版: <https://github.com/thirdwing/ParaBook>

作者和译者尽了最大努力, 但书中错误在所难免

如果对翻译有任何疑问, 请随时通过邮件 kouqiang@mail3.sysu.edu.cn 联系我

关于本书

为什么本书和其它的并行编程书籍不同呢？原因在于我们主要关注在实现层面：

- 这里几乎没有理论内容，诸如 $O()$ 分析、最大理论加速、PRAM、有向无环图 (DAG) 等等。
- 书中使用的都是真实代码。
- 我们使用的都是主流的并行平台，包括 OpenMP、CUDA 和 MPI，而没有使用其它仍处于实验阶段的语言。
- 关于性能的主题——通信延迟、内存/网络连接、负载均衡等，在全书中交叉进行，并且都是在特定平台或应用的层面进行讨论的。
- 相当关注调试技术。

书中使用的主要编程语言 C/C++，但也使用了一些 R 代码。R 已经是最主流的用于数据分析的语言。作为一门脚本语言，R 可以用于快速原型构建。在本书中，我用 R 将可以一些例子表述得远远比用 C/C++ 要简洁，从而使得学生更容易的理解所使用的并行计算原则。出于同样的原因，学生们也可以更容易地编写并行代码，更加集中精力于这些原则之上。另外，R 也有相当丰富的并行库。

我们假设学生在编程方面是有相当经验的，并有包括线性代数在内的数学背景。附录里回顾了本书所需要的数学知识。另一个附录提供了不同系统问题的概述，包括进程调度和虚拟内存等。

需要特别说明的是，书中多数代码没有进行优化。我们主要关注的是技术和语言使用的清晰明了。然而，有很多影响速度的因素值得讨论，比如高速缓存一致性问题、网络延迟、GPU 内存结构等等。

这里展示了你可以如何使用书中的代码：本书使用 L^AT_EX 排版，原始的.tex 文件可以在 <http://heather.cs.ucdavis.edu/~matloff/158/PLN> 下载。请直接下载相关文件（文件名应该足够明确），之后使用一个文本编辑器进行裁剪，从而得到感兴趣的代码。

为了向学生展示研究和教学的相互促进关系，我会时不时引用我的一些研究工作。

如同我的其它开源图书，本书是在不停变动中的。我会继续添加新的主题、新的示例等等，当然也会修补漏洞和改善说明。由于这个原因，所以保存本书最新版本的链接，<http://heather.cs.ucdavis.edu/~matloff/158/PLN/ParProcBook.pdf>，比保存一份拷贝更好。

同样出于这个原因，我非常希望得到反馈。这里我希望感谢 Stuart Ambler、Matt Butner、Stuart Hansen、Bill Hsu、Sameer Khan、Mikel McDaniel、Richard Minner、Lars Seeman、Marc Sosnick 和 Johan Wikström 的评论。特别感谢 Hsu 教授为我提供了高级的 GPU 设备。

各位可能对我另一本关于概览和统计的开源图书感兴趣，可以在 <http://heather.cs.ucdavis.edu/probstatbook> 下载。

本书使用 Creative Commons Attribution-No Derivative Works 3.0 United States License 发行。在美国境外的版权归 Matloff 所有，在保证提供作者和发行信息的情况，这些材料仍可用于教学使用。如果您使用了本书用于教学，我将很高兴您能通知我，这仅仅为了让我知道这些材料正在被使用当中，但这不是必须的。

作者简介

Norm Matloff 博士是加州大学戴维斯分校计算机教授，曾任同校的统计学教授。他曾是硅谷的一名数据库软件开发人员，也曾作为统计咨询师为 Kaiser Permanente Health Plan 等工作过。

Matloff 博士生于 Los Angeles，在 East Los Angeles 和 San Gabriel Valley 长大。他从 UCLA 得到了纯数学的博士学位，研究方向为概率论和统计他在计算机科学和统计学方面发表了大量论文，现在的研究方向是并行处理、统计学和回归方法。

Matloff 博士曾是 UNESCO 下的数据库安全国际委员会，IFIP Working Group 11.3 成员。他是戴维斯分校统计系的创始人之一，并参与了计算机系的建立。他在戴维斯分校被授予 Distinguished Teaching Award and Distinguished Public Service Award。

Matloff 博士是两本书的作者，并编写了大量广泛使用的网络教程，涉及 Linux 和 python 语言。他和 Peter Salzman 博士是《软件调试的艺术》^①一书的作者。Matloff 博士关于 R 语言的《R 语言编程艺术》^②一书已于 2011 年出版。他的新书，*Parallel Computation for Data Science* 会于 2014 年出版。他还写作了很多开源图书，包括 *From Algorithms to Z-Scores: Probabilistic and Statistical Modeling in Computer Science* (<http://heather.cs.ucdavis.edu/probstatbook>)，和《并行机器编程》(<http://heather.cs.ucdavis.edu/~matloff/ParProcBook.pdf>)。

^① 译者注：The Art of Debugging with GDB, DDD, and Eclipse，中文版已由人民邮电出版社出版

^② 译者注：The Art of R Programming，中文版已由机械工业出版社出版

目录

第 1 章 R 并行处理入门	1
1.1 为什么要在本书中用 R 语言?	1
1.2 R 和易并行问题 (Embarrassing Parallel Problems)	1
1.3 一些 R 的并行扩展包	2
1.4 安装和载入这些扩展包	2
1.5 R 中的 snow 扩展包	2
1.5.1 使用	3
1.5.2 示例: 使用 parApply() 进行矩阵向量相乘	3
1.5.3 snow 中的其它函数: clusterApply()、clusterCall() 等	4
1.5.4 示例: 并行求和	6
1.5.5 示例: 对角分块矩阵求逆	8
1.5.6 示例: Mutual Outlink	9
1.5.7 示例: 邻接矩阵变换	11
1.5.8 示例: 设置节点 ID 和集群规模提示	13
1.5.9 关闭集群	14
1.6 multicore 扩展包	14
1.6.1 示例: 使用 multicore 进行邻接矩阵转换	15
1.7 Rdsm	16
1.7.1 示例: 使用 Rdsm 进行对角分块矩阵求逆	17
1.7.2 示例: Web Probe	18
1.7.3 bigmemory 扩展包	19
1.8 R 和 GPU	19
1.8.1 安装	20
1.8.2 gputools 扩展包	20
1.8.3 rgpu 扩展包	20
1.9 通过在 R 中调用 C 进行并行	21
1.9.1 在 R 中调用 C	21
1.9.2 Example: Extracting Subdiagonals of a Matrix	21
1.9.3 在 R 中调用 OpenMPI C 代码	23
1.9.4 在 R 中调用 CUDA 代码	23
1.9.5 示例: Mutual Outlinks	24
1.10 调试 R 程序	25
1.10.1 文本编辑器	25

1.10.2 IDE	25
1.10.3 缺少命令行终端的问题	26
1.10.4 调试 R 所调用的 C 代码	26
1.11 本书中的其它 R 语言示例	26

第 1 章 R 并行处理入门

1.1 为什么要在本书中用 R 语言？

在本书的其它章节里，C/C++ 依然是我们的主要语言，但我们也提供很多 R 语言的示例。为什么要用 R 呢？

- R 是最广泛使用的用于统计分析和数据处理的编程语言。在现今这个大数据时代，人们已经开发了相当数量的用于并行计算的 R 扩展包。特别地，**parallel** 扩展包现在已经是 R 基础包的一部分。
- R 语言的广泛使用，从 Google 设置了其内部的 R 语言规范一事就可见一斑^①。现在 Oracle 也把 R 包含了自己的大数据分析方案中。
- 对于展示各种各样的并行算法，R 非常方便。这点的主要原因在于 R 内置了向量、矩阵和复数类型。

Python 也有很多并行库，比如 **multiprocessing**。关于 Python 的并行话题，我们会在第??章里讨论。

本章的示例会保持尽量简单。但 R 中的并行计算也可以应用到非常庞大而复杂的问题上。在附录??中，有一个 5 分钟的 R 快速入门。阅读时请牢记 R 中 list 结构。

R 中进行并行计算的关键就是——list 结构的操作。许多 R 的并行计算扩展包都非常依赖于 R 中的 list 结构。输入输出的参数和返回值经常都采用 list 的形式。读者可能有兴趣参考一下附录??中的相关内容。

1.2 R 和易并行问题（Embarrassing Parallel Problems）

需要注意的是，R 的并行扩展包一般只能处理易并行问题。正如在 ??节中定义的，这些问题不仅容易并行化，而且信息传递的需求很少^②。如我们所知，一般只有易并行问题会有很好的表现，但在 R 中情况尤其如此，原因如下。

R 语言的函数式编程的本质意味着，任何对一个向量或矩阵的元素的写入操作，比如

```
1 x[3] <- 8
```

都会重写整个向量或矩阵^③。虽然有些例外（随着 R 版本更新，例外可能越来越多），但一般来说我们必须承认 R 中并行的向量和矩阵代码代价很高^④。

对于不易并行的问题，大家应该考虑用 R 调用并行的 C 代码，这点会在 1.9 节中讨论。

^①个人角度来讲，我并不喜欢这些代码规范，我更喜欢我自己的。但从 Google 设置自己的 R 语言规范可以看出他们对 R 的重视程度。

^②后面的要求把很多迭代算法排除在外了，尽管它们很容易并行化。

^③R 中的元素赋值是一个函数调用，上面这个例子的参数分别为 **x**、3 和 8。

^④R 中新引用的类（Reference class）可能会对此有所改变。

1.3 一些 R 的并行扩展包

这里我们列举了一些 R 的并行扩展包：

- Message-passing 或 scatter/gather (??节)：Rmpi、snow、foreach、rmr、Rhipe、multicore^⑤、rzmq
- 内存共享：Rdsm、bigmemory
- GPU：gputools、rgpu

大家可以从 <http://cran.r-project.org/web/views/HighPerformanceComputing.html> 找到更加详尽的列表。

从 2.14 版本开始，R 默认包括了由 snow 和 multicore 构成的 parallel 扩展包。（早期版本可能需要分别下载。）正是因为如此，二者都在范围之内。另外，我们也会讨论 Rdsm/bigmemory 和 gputools。

1.4 安装和载入这些扩展包

安装：

需要注意的是，如果你使用的是 2.14 版或更高版本的 R，你已经安装了 snow 和 multicore。一般来说，除了 rgpu，其它所有扩展包都可以从 R 官方的代码仓库 CRAN (<http://cran.r-project.org>) 下载。这里以 snow 为例：

加入你想把它安装在 /a/b/c/ 目录下。最简单的方法就是使用 R 的函数：

```
1 > install.packages("snow", "/a/b/c/")
```

这会将 snow 安装在 /a/b/c/snow 目录下。

之后你需要将目录 /a/b/c（不是 /a/b/c/snow）加到你的 R 搜索路径中。我推荐大家在自己 home 目录下的 .Rprofile 文件（这是 R 的启动设置文件）中添加这样一行。

```
1 .libPaths("/a/b/c/")
```

在一些情况下，由于所需库的位置原因，你可能需要手动安装一个 CRAN 上的扩展包。这一点请参考下面的 1.8.1 节和 1.8.3 节。

载入一个扩展包：

通过调用 library() 来载入一个扩展包。例如，载入 parallel，可以使用：

```
1 > library(parallel)
```

1.5 R 中的 snow 扩展包

snow 最大的优点在于其简单。其概念和实现都非常简单，能出错的地方不多。因此，它可能是现在使用最广泛的 R 并行包。

snow 扩展包可以直接通过 network socket 运行（由于用户只需要安装 snow，着可能是最常见的用法），也可以运行于 Rmpi（R 的 MPI 接口）、PVM 或 NWS 之上。

^⑤ multicore 扩展包运行于多核内存共享的平台之上，但在读写过程中并不共享数据。

它也可以在一个 scatter/gather 模型（??节）下进行操作。正如 R 中的 `apply()` 函数会将同样的函数作用于一个矩阵的每行上（见下面的示例），`snow` 中的 `parApply()` 会在多台机器上并行地完成类似的操作；不同的机器会操作不同的行。（除了使用多台机器，我们也可以在多核的机器上运行多个 `snow` client。）

1.5.1 使用

在使用

```
1 > library(snow)
```

载入 `snow` 之后，通过调用 `snow` 中的 `makeCluster()` 函数，我们可以设置一个 `snow` 集群。该函数的 `type` 参数用于选择网络平台，诸如“MPI”或“SOCK”。后者用于将 `snow` 运行于其自己创建的 TCP/IP sockets 之上，而不是使用 MPI。

在这个例子里，我在名为 `pc48` 和 `pc49` 的电脑上使用“SOCK”选择，以这种方式设置集群^⑥：

```
1 > cls <- makeCluster(type="SOCK",c("pc48","pc49"))
```

需要注意的是上面的 R 代码在名为 `pc48` 和 `pc49` 的机器上设置了工作节点；这和管理节点相区别，管理节点运行于执行 R 代码的机器上。

如果你想把工作节点和管理节点同时运行在同一台机器上（特别是在一台多核的机器上），需要使用 `localhost` 作为机器名。

还有其它很多可选的参数。一个你可能觉得非常有用的是 `outfile`，它会把调用的结果记录在名为 `outfile` 的文件里。这在调用失败进行 debug 时非常有用。

1.5.2 示例：使用 `parApply()` 进行矩阵向量相乘

为了介绍 `snow`，让我们考虑一个简单的矩阵向量相乘的简单示例。我是指一个测试矩阵如下：

```
1 > a <- matrix(c(1:12),nrow=6)
2 > a
3      [,1] [,2]
4 [1,]  1  7
5 [2,]  2  8
6 [3,]  3  9
7 [4,]  4 10
8 [5,]  5 11
9 [6,]  6 12
```

我们会将向量 $(1,1)^T$ （T 这里表示转置）和矩阵相乘。在这个简单的示例，我们当然可以直接完成：

```
1 > a %*% c(1,1)
2      [,1]
3 [1,]  8
```

^⑥ 如果你使用的是一个文件共享系统的电脑集群，尽量保证 R 的安装路径一致，以避免问题。


```

4 [2,] 10
5 [3,] 12
6 [4,] 14
7 [5,] 16
8 [6,] 18

```

但是让我们看看如何使用 R 的 `apply()` 来完成它。尽管这仍是顺序执行，但这为我们扩展到并行计算提供了便利。

R 的 `apply()` 函数调用一个用户定义的标量函数作用于用户指定的矩阵的每一行（或每一列）。为了将 `apply()` 用于这里的矩阵向量相乘问题，我们定义一个点积的函数：

```
1 > dot <- function(x,y) {return(x%*%y)}
```

现在调用 `apply()`：

```

1 > apply(a,1,dot,c(1,1))
2 [1] 8 10 12 14 16 18

```

这个调用将函数 `dot()` 作用于矩阵 `a` 的每一行（这个可以从 1 看出，2 意味着每一列）；每一行都将作为 `dot()` 的第一个参数，而 `c(1,1)` 会作为第二个参数。换言之，`dot()` 的第一次调用就是

```
1 dot(c(1,7),c(1,1))
```

`snow` 中的 `parApply()` 函数将 `apply()` 扩展到并行计算。我们把它用于将我们的矩阵相乘问题并行化，运行在我们名为 `cls` 的集群之上：

```

1 > parApply(cls,a,1,dot,c(1,1))
2 [1] 8 10 12 14 16 18

```

`parApply()` 所作的就是将矩阵每一行发送给每一个节点，同时发送的还由函数 `dot()` 和参数 `c(1,1)`。每个节点将 `dot()` 作用到接收的行上，之后将结果返回给管理节点。

R 的 `apply()` 函数一般只用于变量值的情形，这意味着 `apply(m,i,f)` 调用中的函数 `f()` 的返回值是标量。如果 `f()` 的返回值是向量值，那返回的会是一个矩阵而不是一个向量，矩阵里的每一列是 `f()` 作用于 `m` 的一列或一行的结果。`parApply()` 也同样如此。

1.5.3 snow 中的其它函数：clusterApply()、clusterCall() 等

上一节，我们介绍了 `parApply()` 函数。它可以这样调用

- `parApply()`:

```
1 parApply(cls,m,DIM,f,...)}
```

这个调用会把矩阵 `m` 的每一行分配到 `cls` 的各个工作节点，之后函数 `f()` 会被作用到每一行，省略号在这里表示可选参数。参数 `DIM` 为 1 时表示行操作，2 表示列操作。

返回值是一个向量（也可能是个矩阵，如上所述）。

`snow` 最大的有点在于其简单，因此并没有很多复杂的函数，但当然不止 `parApply()` 一个。这里列举了一些：

- `clusterApply()`:

这个函数可能是 `snow` 中被使用最频繁的函数。

```
1 clusterApply(cls, individualargs, f, ...)
```

这会使 `f()` 在 `cls` 中的每个节点上运行。这里的 `individualargs` 是一个 R 列表（如果是个向量，会被转换成列表）。当 `f()` 在集群中的节点 `i` 上被调用时，其参数如下所述：第一个参数是 `individualargs` 的第 `i` 个元素，或者说是 `individualargs[[i]]`；如果在调用时，是用了省略号所代表的（可选）参数，它们会作为第二、第三或更多的参数传递给 `f()`。

如果 `individualargs` 的元素数量大于集群中的节点数，那么 `cls` 会被循环使用（可以把它作为一个向量对待），所以多数或全部节点会在不止一个 `individualargs` 元素上调用 `f()`。返回值是一个 R 列表，其中第 `i` 个元素是 `f()` 作用于 `individualargs` 中第 `i` 个元素的结果。所以说，`individualargs` 列表又需要拆分并行计算的工作构成。

- **clusterApplyLB():**

这是 `clusterApply()` 的负载均衡模式，目的在于解决我们在第??章中提到的性能问题。

为了解释 `clusterApply()` 的两者形式的区别，假设我们的集群由 10 个节点，而我们有 25 个需要执行的任务（或者说 `individualargs` 的长度是 25）。如果使用 `clusterApply()`，会发生下列这些：

- 前 10 个任务会被分配给工作节点，每个节点一个任务。
- 管理节点会等这 10 个任务完成，之后再分配另外 10 个。
- 管理节点会等这 10 个任务完成，之后在分配剩下的 5 个。
- 管理节点会等这 5 个任务完成，之后返回 25 个结果。

而是用 `clusterApplyLB()` 时，会按照下面这种方式执行：

- 前 10 个任务会被分配给工作节点，每个节点一个任务。
- 当由节点任务结束时，管理节点会马上行动，将第 11 个任务分配给这个节点，即使其它节点的任务还没完成。
- 管理节点会继续照此工作，一旦一个节点任务完成，就会分配新的任务，知道所有任务完成。
- 管理节点最后会返回 25 个结果。

用第??章和 OpenMP 一章中的??节的说法，`clusterApply()` 使用了一种静态的调度策略，而 `clusterApplyLB()` 使用了一种动态策略；其中 chunk size 为 1。

- **clusterCall():**

函数 `clusterCall(cls, f, ...)` 将函数 `f()` 和省略号所代表的参数（如果有的话），发送到每个工作节点。在每个节点上，`f()` 会使用这些参数求值。返回值是一个 R 列表，第 `i` 个元素是第 `i` 个节点的计算结果。（一眼看上去，似乎每个节点都会返回同样的结果，但 `f()` 会使用每个节点特定的参数，从而返回不同的结果。）

- **clusterExport():**

函数 `clusterExport(cls, varlist)` 会将名字出现在字符向量 `varlist` 中的变量拷贝到 `cls` 中的各个节点。你可以使用这个函数来避免从管理节点到工作节点开销巨大的数据传输。使用这个函数，你可以只传输数据集一次；通过在相应的变量上使用 `clusterExport()`，之后在工作节点上将其作为全局变量使用。同样地，返回值仍是个 R 列表，第 `i` 个元素是集群中第 `i` 个节点的计算结果。

默认情况下，被传输到工作节点的变量在管理节点上必须是全局变量。

需要特别注意的是，一旦你传输了一个变量，比如 `x`，从管理节点到各个工作节点上，各个拷贝和工作节点上的变量就是独立的了（各个拷贝之间也是相互独立的）。如果其中一个拷贝改变了，在其他拷贝中不会反应这些变化。

- `clusterEvalQ()`:

函数 `clusterEvalQ(cls,expression)` 会在 `cls` 的各个节点上运行 `expression`。

1.5.4 示例：并行求和

现在让我们再看一个示例，我们用 `snow` 来进行并行求和。先从一个很简单的版本开始，之后再考虑复杂的版本。

```
1 parsum <- function(cls,x) {
2   # 在节点上分配 x 的索引（实际上没有传输任何东西）
3   xparts <- clusterSplit(cls,x)
4   # 现在传输到节点上，并进行求和
5   tmp <- clusterApply(cls,xparts,sum)
6   # 现在将各个单独的加和合并得到结果
7   tot <- 0
8   for (i in 1:length(tmp)) tot <- tot + tmp[[i]]
9   return(tot)
10 }
```

现在我们在一个有两个共走节点的集群 `cls` 上进行测试：

```
1 > x
2 [1] 1 2 3 4 5 6 5 12 13
3 > parsum1(cls,x)
4 [1] 51
```

结果不错。现在我们来想一下，这是如何完成的？

最基本的想法就是讲我们的向量分块，之后分配给工作节点。每个工作节点会把所分配的小块求和，再把结果返回给管理节点。管理节点会把这些结果求和，返回我们想要的求和的最终结果。

为了将我们的向量 `x` 分块并发给各个节点，我先来看 `snow` 中的函数 `clusterSplit()`。这个函数的输入是一个 R 向量，之后将其分块，分块的数量和工作节点数相同。

例如，在上面的两个工作节点的集群上，我们得到：

```
1 > xparts <- clusterSplit(cls,x)
2 > xparts
3 [[1]]
4 [1] 1 2 3 4
5
6 [[2]]
7 [1] 5 6 5 12 13
```

非常肯定的是，我们的列表 `xparts` 有在其一个元素中有 `x` 的一块，而另一个元素中有 `x` 的另一块。之后这两块被传输到两个工作节点上：

```
1 > tmp <- clusterApply(cls,xparts,sum)
2 > tmp
3 [[1]]
4 [1] 10
5
6 [[2]]
7 [1] 41
```

同样像 `snow` 中的其他函数一样，`clusterApply()` 会以列表的形式返回结果。这里我们将结果赋值给了 `tmp`。其内容如下

```
1 > tmp
2 [[1]]
3 [1] 10
4
5 [[2]]
6 [1] 41
```

也就是 `x` 每一小块的和。

为了得到最后的结果，我们不能简单地在 `tmp` 上使用 R 中 `sum()` 函数：

```
1 > sum(tmp)
2 Error in sum(tmp) : invalid 'type' (list) of argument
```

这是因为 `sum()` 接受的是向量，而不是列表。所以我们自己写一个循环来把结果加起来：

```
1 tot <- 0
2 for (i in 1:length(tmp)) tot <- tot + tmp[[i]]
```

需要注意的一点是，我们使用 `[[]]` 来获取列表中的元素。

我可以通过调用 R 中的 `Reduce()` 函数来取代上面的循环，从而对代码进行优化。`Reduce()` 很像??节和??节中的 reduction 操作。（注意，这里是个串行操作，不是并行。）一般以 `Reduce(f,y)` 这种形式使用，它对函数 `f()` 和列表 `y` 进行如下操作

```
1 z <- y[1]
2 for (i in 2:length(y)) z <- f(z,y[i])
```

使用 `Reduce()` 可以使代码更紧凑可读，一些情况下还会提高执行效率（我们这里只有很少的项目进行相加，暂时还不用考虑效率）。而且，`Reduce()` 会将 `tmp` 从一个列表转换为向量，这就解决了我们直接对 `tmp` 使用 `sum()` 时的的问题。

下面是新的代码：

```
1 parsum <- function(cls,x) {
2   xparts <- clusterSplit(cls,x)
3   tmp <- clusterApply(cls,xparts,sum)
```

```

4   Reduce(sum,tmp) # implicit return()
5 }

```

需要说明的是，在 R 中，如果没有显式的 `return()` 语句，那最后求得值会被作为返回值，这里是 `Reduce()` 的计算结果。

`Reduce()` 是一个非常便利的函数，特别是在和 `snow` 一起使用时。这里有一个我们把多个矩阵进行合并的示例：

```

1 > Reduce(rbind,list(matrix(5:8,nrow=2),3:4,c(-1,1)))
2   [,1] [,2]
3 [1,] 5 7
4 [2,] 6 8
5 [3,] 3 4
6 [4,] -1 1

```

`rbind()` 函数只有两个参数，在这里我们三个。通过使用 `Reduce()` 可以解决这个问题。

1.5.5 示例：对角分块矩阵求逆

假设我们有一个对角分块矩阵，比如

$$\begin{pmatrix} 1 & 2 & 0 & 0 \\ 3 & 4 & 0 & 0 \\ 0 & 0 & 8 & 1 \\ 0 & 0 & 1 & 5 \end{pmatrix}$$

我们想对其求逆。这是个易并行问题：假如我们有两个处理器，我们可以很简单地让其中之一对第一个 2×2 子矩阵求逆，让另一个对第二个 2×2 子矩阵求逆，之后我们将两个逆矩阵放回原来的位置。

通讯的开销在这里不是很大，一个 $n \times n$ 矩阵求逆的时间复杂度为 $O(n^3)$ ，而通讯只有 $O(n^2)$ 。

现在我们讨论一下用于分块对角矩阵求逆的 `snow` 代码。

```

1 # invert a block diagonal matrix m, whose sizes are given in szs;
2 # return value is the inverted matrix
3 bdiaginv <- function(cls,m,szs) {
4   nb <- length(szs) # number of blocks
5   dgs <- list() # will form args for clusterApply()
6   rownums <- getrng(szs)
7   for (i in 1:nb) {
8     rng <- rownums[i,1]:rownums[i,2]
9     dgs[[i]] <- m[rng,rng]
10  }
11  invs <- clusterApply(cls,dgs,solve)
12  for (i in 1:nb) {
13    rng <- rownums[i,1]:rownums[i,2]
14    m[rng,rng] <- invs[[i]]
15  }

```

```

16     m
17 }
18
19 # find row number ranges for the blocks, returned in a # 2-column
20 # matrix; blkzs = block sizes
21 getrng <- function(blkzs) {
22     col2 <- cumsum(blkzs) # cumulative sums function
23     col1 <- col2 - (blkzs-1)
24     cbind(col1,col2) # column bind
25 }

```

我们来检测一下：

```

1 > m
2     [,1] [,2] [,3] [,4] [,5]
3 [1,] 1 2 0 0 0
4 [2,] 7 8 0 0 0
5 [3,] 0 0 1 2 3
6 [4,] 0 0 2 4 5
7 [5,] 0 0 1 1 1
8 > bdiaginv(cls,m,c(2,3))
9     [,1] [,2] [,3] [,4] [,5]
10 [1,] -1.333333 0.3333333 0 0 0
11 [2,] 1.166667 -0.1666667 0 0 0
12 [3,] 0.000000 0.0000000 1 -1 2
13 [4,] 0.000000 0.0000000 -3 2 -1
14 [5,] 0.000000 0.0000000 2 -1 0

```

这里的 `szs` 参数，包含了分块的大小。由于我们只有一个 2×2 和一个 3×3 的块，分块的大小就是 2 和 3，因为在函数调用里使用 `c(2,3)`。

这里 `clusterApply()` 的使用和早先的例子很相似。代码中值得注意的地方是我们需要保存每一块在大矩阵中的位置。最后我们写了一个 `getrng()` 函数，用于返回不同块的起始和结束的行数。我们通过使用这个函数来设置 `clusterApply()` 的 `dg` 参数：

```

1 for (i in 1:nb) {
2     rng <- rownums[i,1]:rownums[i,2]
3     dgs[[i]] <- m[rng,rng]

```

大家要记得表达式 `m[rng,rng]` 会提取 `m` 的行和列出来，在这里就是第 `i` 块。

1.5.6 示例：Mutual Outlink

让我们考虑??节中的例子。我们有一个网络，比如 web 链接。对于其中的两个节点，比如两个网站，我可能对其 mutual outlink 感兴趣，也就是两个网站共同的对外链接。

下面的 `snow` 代码会计算整个网络中任意一对节点 mutual outlink 的均值。

```

1 # snow version of mutual links problem
2
3 library(snow)
4
5 mtl <- function(ichunks,m) {
6   n <- ncol(m)
7   matches <- 0
8   for (i in ichunks) {
9     if (i < n) {
10      rowi <- m[i,]
11      matches <- matches +
12        sum(m[(i+1):n,] %*% as.vector(rowi))
13    }
14  }
15  matches
16 }
17
18 # returns the mean number of mutual outlinks in m, computing on the
19 # cluster cls
20 mutlinks <- function(cls,m) {
21   n <- nrow(m)
22   nc <- length(cls)
23   # determine which worker gets which chunk of i
24   options(warn=-1)
25   ichunks <- split(1:n,1:nc)
26   options(warn=0)
27   counts <- clusterApply(cls,ichunks,mtl,m)
28   do.call(sum,counts) / (n*(n-1)/2)
29 }

```

对于 **m** 中的每一行，我们会计算其下面每一行中的 mutual link。为了在工作节点之间分配工作，我们可以如下使用 **clusterSplit()**

```
1 clusterSplit(cls,1:nrow(m))
```

但这会有一个在??节中讨论过的不均衡问题。比如我们有两个工作节点和 100 行。如果我们像上一节一样使用 **clusterSplit()**，第一个节点进行的比较工作会远比第二个节点多。

一个解决方案是在调用 **clusterSplit()** 之前，将行号随机打乱。另一方法，也是我们上面的代码中使用的，就是用 R 的 **split()** 函数。

那 **split()** 是做什么的？它根据第二个参数中设置的“类别”，将第一个参数进行分块处理。我们来看这个示例：

```

1 > split(2:5,c('a','b'))
2 $a
3 [1] 2 4

```

```

4
5 $b
6 [1] 3 5

```

这里的种类是 a 和 b。`split()` 函数要求第二个参数和第一个参数长度相同，所以首先会对第二个参数进行“循环”处理成 a,b,a,b,a。之后会将 2 和 4 放入类别 a，将 3 和 5 放入类别 b。`split()` 函数最后会返回一个相应的列表。

现在再回到上面的 `snow` 示例，我们仍然假设在两个工作节点中分配 100×100 的矩阵 `m`，代码

```

1 nc <- length(cls)
2 ichunks <- split(1:n,1:nc)

```

会生成一个由两部分构成的列表，第一部分由奇数行构成，第二部分由偶数行构成。之后我们再使用

```
1 counts <- clusterApply(cls,ichunks,mtl,m)
```

就可以在两个工作节点间实现一个负载均衡了。

注意这个调用需要将 `m` 作为一个参数（作为函数 `mtl()` 的参数）。否则工作节点将没有可供使用的 `m`。另一个选择是使用 `clusterExport()` 来将 `m` 发送到工作节点，之后作为一个全局变量供 `mtl()` 使用。

另外，调用 `options()` 是为了让 R 在我们做“循环”时不发出警告。一般我们并不这么做，但这里为了使用 `split()` 的需要。

之后，为了得到输出的列表中各个元素的总和，我们可以再次使用 `Reduce()`，但由于 R 的多样性，我们也可以使用 `do.call()` 函数。这个函数的动能正如其函数名暗示的：它会把列表 `counts` 中的每个元素抽出，之后作为参数传递给 `sum()`！（一般来说，当我们需要调用一个特定的函数，但其参数的数目直到运行时才可以确定时，`do.call()` 是非常有用的。）

正如前面所说的，除了使用 `split()`，我们可以将行数随机打乱：

```
1 tmp <- clusterSplit(cls,order(runif(nrow(m))))
```

这会为每一行产生一个 (0,1) 之间的随机数，之后按此排序。比如说，如果第三个随机数是第 20 小的，第三个元素在 `order()` 的输出中会是 20。这可以找到矩阵 `m` 行号的一个随机排列。

1.5.7 示例：邻接矩阵变换

这是??节中代码的 `snow` 版本。回忆一下，问题如下：

假如我们有一个图的邻接矩阵

$$\begin{pmatrix} 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 1 \\ 0 & 1 & 0 & 1 \\ 1 & 1 & 1 & 0 \end{pmatrix} \quad (1.1)$$

其中行号和列号从 0 开始，而不是 1。我们想要将其转换为一个两列的矩阵用来展示连接数，如下

所示

$$\begin{pmatrix} 0 & 1 \\ 1 & 0 \\ 1 & 3 \\ 2 & 1 \\ 2 & 3 \\ 3 & 0 \\ 3 & 1 \\ 3 & 2 \end{pmatrix} \quad (1.2)$$

比如说，在上面的邻接矩阵中，最右边的第二行有一个 1，这意味着在顶点 1 和顶点 3 直接存在一条边。这个在转换后的矩阵中以 (1,3) 表示。

下面是在 **snow** 中进行该计算的代码：

```

1  tg <- function(cls,m) {
2    n <- nrow(m)
3    rowschunks <- clusterSplit(cls,1:n) # make chunks of row numbers
4    m1 <- cbind(1:n,m) # prepend col of row numbers to m
5    # now make the chunks of rows themselves
6    tmp <- lapply(rowschunks,function(rchunk) m1[rchunk,])
7    # launch the computation
8    tmp <- clusterApply(cls,tmp,tgonchunk)
9    do.call(rbind,tmp) # combine into one large matrix
10 }
11
12 # a worker works on a chunk of rows
13 tgonchunk <- function(rows) {
14   # note: matrix space allocation not efficient
15   mat <- NULL
16   nc <- ncol(rows)
17   for (i in 1:nrow(rows)) {
18     row <- rows[i,]
19     rownum <- row[1]
20     for (j in 2:nc) {
21       if (row[j] == 1) {
22         if (is.null(mat)) {
23           mat <- matrix(c(rownum,j-1),ncol=2)
24         } else
25           mat <- rbind(mat,c(rownum,j-1))
26       }
27     }
28   }
29   return(mat)
30 }
```

这里有什么新东西么？首先，由于我们需要最后的输出按字典序排列，我们需要保存原有每行的索引。所以我们需要在 `m` 中多添加一行：

```
1 m1 <- cbind(1:n,m) # prepend col of row numbers to m
```

其次，注意 `lapply()` 函数的使用。正如 `apply()` 会把一个特定的函数作用到矩阵的每一行（或每一列）上，`lapply()` 会把一个特定函数作用到列表的每一个元素上。输出结果依然是个列表。

在我们这里的例子中，我们需要将 `m` 按行分块传递给 `clusterApply()`，但后者要求我们必须传递一个列表。我们可以通过一个 `for` 来完成，一个一个地将分块添加进列表，但使用 `lapply()` 可以更加紧凑。

在最后，管理节点会接收新矩阵的很多部分，这些必须被整合起来。使用 `rbind()` 函数是很自然的想法，但我们仍然需要克服各个部分是 R 列表的问题。尽管 `Reduce()` 也可以完成，但用 `do.call()` 会更趁手。

需要注意的是，尽管在上一段中说使用 `rbind()` 是很自然的，但效率很低。这是因为 `rbind()` 会重新分配一个新的矩阵空间，这是个很浪费时间的操作。先分配 50 行空间，之后在构建矩阵时进行填充会是更好的选择。无论什么时候，我们用完了矩阵，我们都可以构建一个新的矩阵，之后把所有矩阵作为一个列表返回。

1.5.8 示例：设置节点 ID 和集群规模提示

让我们回忆一下，在 OpenMP 中有两个函数，`omp_get_thread_num()` 和 `omp_get_num_threads()`，分别用来报告一个线程的 ID 和线程总数。在 MPI 中，对应的函数是 `MPI_Comm_rank()` 和 `MPI_Comm_size()`。在 `snow` 中如果能有这样的函数（或功能），将是非常好的事情。这里的代码就是用于这个目的：

```
1 # sets a list myinfo as a global variable in the worker nodes in the
2 # cluster cls, with myinfo$id being the ID number of the worker and
3 # myinfo$nwkrks being the number of workers in the cluster; called from
4 # the manager node
5 setmyinfo <- function(cls) {
6   setmyinfo <- function(i,n) {
7     myinfo <-<- list(id = i, nwkrks = n)
8   }
9   ncls <- length(cls)
10  clusterApply(cls,1:ncls,setmyinfo,ncls)
11 }
```

是的，R 允许在函数中定义函数。顺便请注意超级赋值符 `<-<` 的使用，这个用了在全局层面进行赋值操作。

调用这个函数后，任何在一个工作节点上运行的代码代码都可以决定其节点 ID，比如在下面这样的代码中

```
1 if (myinfo$id == 1) ...
```

或者，我们也可以从管理节点传输代码到工作节点执行：

```
1 > setmyinfo(cls)
2 [[1]]
3 [[1]]$id
4 [1] 1
5
6 [[1]]$nwrkrs
7 [1] 2
8
9
10 [[2]]
11 [[2]]$id
12 [1] 2
13
14 [[2]]$nwrkrs
15 [1] 2
16
17 > clusterEvalQ(cls,myinfo$id)
18 [[1]]
19 [1] 1
20
21 [[2]]
22 [1] 2
```

第一个例子，由于 `clusterApply()` 有返回值，都会被打印出来。第二个例子中，调用

```
1 clusterEvalQ(cls,myinfo$id)
```

会使每个工作节点对表达式 `myinfo$id` 进行求值；之后 `clusterEvalQ()` 返回在每个节点上的执行结果。

1.5.9 关闭集群

退出 R 之前，不要忘记使用 `stopCluster(clustername)` 来关闭集群。

1.6 multicore 扩展包

正如名字所暗示的，**multicore** 扩展包就是用于使用多核设备的计算能力的。这可能有点奇怪：既然 **snow** 既可以用于一个（物理）集群，也可以用于一个多核设备，而 **multicore** 只能在后者上使用，那用 **multicore** 的优势哪儿？答案是性能上的提高，这个在后面会解释。

这个扩展包的主函数是 `mclapply()`，其语法和 **snow** 中的 `clusterApply()` 很类似，也很类似地把任务分配给各个工作节点。

这里所说的工作节点，指的是同一台机器上的不同处理器。比如说，在一个四核的机器上运行 **multicore**，调用 `mclapply()` 会（默认）在你的机器上调用 4 个 R，并行地进行你的运算工作。其中每个 R 调用都使用和调用前的 R 一样的变量设置。因此所有的变量最初（注意这个修饰语）

是共享的，而不需要程序员采取特别的措施来将变量从管理节点分配到工作节点，这和 **snow** 相当不同。

这一切都是由 **mclapply()** 调用你操作系统中的 **fork()** 函数来完成的。（所以这仅限于类 Unix 系统，比如 Linux 和 MacOS。）这个 fork 过程是由 R 自己完成的，每个工作节点一个新的拷贝。

因此分配到 R 拷贝的工作节点会共享所有在 fork 发生时存在的变量（也包括你调用 **mclapply()** 时的局部变量）。所以你的代码不需要将这些变量拷贝到工作节点，工作节点会自动获取它们。但需要注意的是，这些变量只是在最初的时候是共享的，对其中一个拷贝的修改不会再其它拷贝中有所反应（包括最初的那个）。

从管理节点到工作节点的初始值拷贝是基于 **copy-on-write** 的，这意味着直到一个节点尝试获取数据，这份数据才会被拷贝过去。这个粒度（granularity）是在虚拟内存页（virtual memory page）层面（??节）上的。同样，这由操作系统处理，不是 R。

由于这个物理拷贝最终还是会由操作系统完成，所以 **multicore** 相对 **snow** 并没有很多人所想的那么有优势。然而，这可能在处理延迟方面由一定优势（??节）。有些情况下，不需要所有节点同时获取变量，所以可能一个节点在拷贝变量，而其余的在进行计算。

还需要注意一点，**snow** 中，一个集群被设置好，会在每个 **snow** 函数调用中重复使用，而 **multicore** 与此不同，每一个 R 进程都在一个 **multicore** 函数被调用时从头开始。

1.6.1 示例：使用 multicore 进行邻接矩阵转换

和1.5.7节中的示例一样，而且实际上下面的 **tgonchunk()** 函数就是我们前面 **snow** 代码的修改版。

```
1 mclapply(starts,tgonchunk,m1,chunksize,mc.cores=ncores)
```

这个调用会将 **tgonchunk()** 函数作用于 **starts** 向量的每一个元素上（首先会被转换成 R 列表），其中 **m1** 和 **chunksize** 作为 **mclapply()** 的参数使用。

```
1 # transgraph problem, R multicore version
2
3 # arguments:
4 # m: the input matrix
5 # ncores: desired number of cores to use
6 tgmcm <- function(m,ncores) {
7   n <- nrow(m)
8   chunksize <- floor(n/ncores)
9   starts <- seq(1,n,chunksize)
10  m1 <- cbind(1:n,m) # prepend col of row numbers to m
11  tmp <- mclapply(starts,tgonchunk,m1,chunksize,mc.cores=ncores)
12  do.call(rbind,tmp)
13 }
14
15 # a worker works on a chunk of rows
16 tgonchunk <- function(start,m1,chunksize) {
17   # note: matrix space allocation not efficient
```

```

18   outmat <- NULL
19   end <- start + chunksize - 1
20   nrm <- nrow(m1)
21   if (end > nrm) end <- nrm
22   ncm <- ncol(m1)
23   for (i in start:end) {
24     rownum <- m1[i,1]
25     for (j in 2:ncm) {
26       if (m1[i,j] == 1) {
27         if (is.null(outmat)) {
28           outmat <- matrix(c(rownum,j-1),ncol=2)
29         } else
30           outmat <- rbind(outmat,c(rownum,j-1))
31       }
32     }
33   }
34   return(outmat)
35 }

```

1.7 Rdsm

无论你在一个 NOW 网络还是一个多核机器上，我的 **Rdsm** 扩展包都可以作为一个多线程来使用。这是我在 2002 年开发的一个类似的 Perl 扩展包，PerlDSM^⑦的扩展。**Rdsm** 的主要优势在于：

- 使用了一个内存共享的编程模型，正如在??节中所述，在并行编程社区中，一般认为这优于信息传递模型。
- 可以充分使用 R 的调试工具。

Rdsm 给了 R 程序员一个内存共享的视角，但实际上这些对象并没有共享。对象被储存在一个服务器上，通过网络端口获取^⑧，从而使 R 程序员即使在 NOW 网络上也可以有一个类似多线程的视角。这里没有管理/工作节点的结构，所有的 R 进场都执行相同的代码。

Rdsm 中的共享对象，可以是 **dsmv** 和 **dsmm** 类中的数值向量或矩阵，也可以是 **dsml** 类中的 R 列表。为了效率，向量和矩阵与服务器的通讯是二进制的形式进行的，而列表进行了序列化。还有一个内置变量 **myinfo** 用于获取每一个进程的 ID 和进程总数，这和 **Rmpi** 中的 **mpi.comm.rank()** 和 **mpi.comm.size()** 返回的信息一样。

Rdsm 同样可以使用上面提到的 **install.packages()** 进行安装。**Rdsm** 提供了内置文档，不过最好还是要通读 **examples** 文件夹下的 **MatMul.R** 代码。里面提供了大量注释，希望可以作为这个扩展包的一个入门。

^⑦N. Matloff, PerlDSM: A Distributed Shared Memory System for Perl, *Proceedings of PDPTA 2002*, 2002, 63-68.

^⑧**Rdsm** 也可以在 **bigmemory** 扩展包中使用，见1.7.3节。

1.7.1 示例：使用 Rdsm 进行对角分块矩阵求逆

现在让我们来看如何将1.5.5节中的对角分块矩阵求逆使用 **Rdsm** 处理。

```

1 # invert a block diagonal matrix m, whose sizes are given in szs; here m
2 # is either an Rdsm or bigmemory shared variable; no return
3 # value--inversion is done in-place; it is assumed that there is one
4 # thread for each block
5
6 bdiaginv <- function(bd,szs) {
7   # get number of rows of bd
8   nrdb <- if(class(bd) == "big.matrix") dim(bd)[1] else bd$size[1]
9   rownums <- getrng(nrdb,szs)
10  myid <- myinfo$myid
11  rng <- rownums[myid,1]:rownums[myid,2]
12  bd[rng,rng] <- solve(bd[rng,rng])
13  barr() # barrier
14 }
15
16 # find row number ranges for the blocks, returned in a 2-column matrix;
17 # matsz = number of rows in matrix, blkszs = block sizes
18 getrng <- function(matsz, blkszs) {
19   nb <- length(blkszs)
20   rwnms <- matrix(nrow=nb,ncol=2)
21   for (i in 1:nb) {
22     # i-th block will be in rows (and cols) i1:i2
23     i1 <- if (i==1) 1 else i2 + 1
24     i2 <- if (i == nb) matsz else i1 + blkszs[i] - 1
25     rwnms[i,] <- c(i1,i2)
26   }
27   rwnms
28 }

```

相较于 **snow** 中的 11 行代码，这里主要的并行工作由这 4 行完成：

```

1 myid <- myinfo$myid
2 rng <- rownums[myid,1]:rownums[myid,2]
3 bd[rng,rng] <- solve(bd[rng,rng])
4 barr() # barrier

```

这也展示了内存共享编程模型相对信息传递模型的优势。

1.7.2 示例: Web Probe

In the general programming community, one major class of applications, even on a serial platform, is parallel I/O. Since each I/O operation may take a long time (by CPU standards), it makes sense to do them in parallel if possible. **Rdsm** facilitates doing this in R.

The example below repeatedly cycles through a large list of Web sites, taking measurements on the time to access each one. The data are stored in a shared variable **accesstimes**; the **n** most recent access times are stored. Each **Rdsm** process works on one Web site at a time.

An unusual feature here is that one of the processes immediately exits, returning to the R interactive command line. This allows the user to monitor the data that is being collected. Remember, the shared variables are still accessible to that process. Thus while the other processes are continually adding data to **accesstimes** (and deleted one item for each one added), the user can give commands to the exited process to analyze the data, say with histograms, as the collection progresses.

Note the use of lock/unlock operations here, with the **Rdsm** variables of the same names.

```

1 # if the variable accesstimes is length n, then the Rdsm vector
2 # accesstimes stores the n most recent probed access times, with element
3 # i being the i-th oldest
4
5 # arguments:
6 # sitefile: IPs, one Web site per line
7 # ww: window width, desired length of accesstimes
8 webprobe <- function(sitefile,ww) {
9   # create shared variables
10   cnewdsm("accesstimes","dsmv","double",rep(0,ww))
11   cnewdsm("naccesstimes","dsmv","double",0)
12   barr() # Rdsm barrier
13   # last thread is intended simply to provide access to humans, who
14   # can do analyses on the data, typing commands, so have it exit this
15   # function and return to the R command prompt
16   # built-in R list myinfo has components to give thread ID number and
17   # overall number of threads
18   if (myinfo$myid == myinfo$nclnt) {
19     print("back to R now")
20     return()
21   } else { # the other processes continually probe the Web:
22     sites <- scan(sitefile,what="") # read from URL file
23     nsites <- length(sites)
24     repeat {
25       # choose random site to probe
26       site <- sites[sample(1:nsites,1)]
27       # now probe it, recording the access time

```

```

28     acc <- system.time(system(paste("wget --spider -q",site)))[3]
29     # add to accesstimes, in sliding-window fashion
30     lock("acclock")
31     if (nacesstimes[1] < ww) {
32         nacesstimes[1] <- nacesstimes[1] + 1
33         accesstimes[nacesstimes[1]] <- acc
34     } else {
35         # out with the oldest, in with the newest
36         newvec <- c(accesstimes[-1],acc)
37         accesstimes[] <- newvec
38     }
39     unlock("acclock")
40 }
41 }
42 }

```

1.7.3 bigmemory 扩展包

Jay Emerson 和 Mike Kane 在我开发 **Rdsm** 的同时，开发了 **bigmemory** 扩展包；而我们互相都不知道这一点。

The **bigmemory** package is not intended to provide a threads environment. Instead, it is used to deal with a hard limit R has: No R object can be larger than $2^{31} - 1$ bytes. This holds even if you have a 64-bit machine with lots of RAM. The **bigmemory** package solves the problem on a multicore machine, by making use of operating system calls to set up shared memory between processes.^⑨

In principle, **bigmemory** could be used for threading, but the package includes no infrastructure for this. However, one can use **Rdsm** in conjunction with **bigmemory**, an advantage since the latter is very efficient.

Using **bigmemory** variables in **Rdsm** is quite simple: Instead of calling **cnewdsm()** to create a shared variable, call **newbm()**.

1.8 R 和 GPU

The blinding speed of GPUs (for certain problems) is sure to of interest to more and more R users in the coming years.

As of today, the main vehicle for writing GPU code is CUDA, on NVIDIA graphics cards. CUDA is a slight extension of C.

You may need to write your own CUDA code, in which case you need to use the methods of Section 1.9. But in many cases you can get what you need in ready-made form, via the two main packages for GPU programming with R, **gputools** and **rgpu**. Both deal mainly with linear algebra operations. The remainder of this section will deal with these packages.

^⑨It can also be used on distributed systems, by exploiting OS services to map memory to files.

1.8.1 安装

Note that, due to issues involving linking to the CUDA libraries, in the cases of these two packages, you probably will *not* be able to install them by merely calling `install.packages()`. The alternative I recommend works as follows:

- 下载 `.tar.gz` 格式的扩展包。
- 将扩展包解压缩，我们把产生的文件夹叫做 `x`。
- 假设你想把它安装到 `/a/b/c`。
- 对 `x` 中的文件进行修改。
- 之后运行 R CMD INSTALL -l /a/b/c x。

细节会在后面的章节中讨论。

1.8.2 gputools 扩展包

In installing `gputools`, I downloaded the source from the CRAN R repository site, and unpacked as above. I then removed the subcommand

```
-gencode arch=compute_20,code=sm_20
```

from the file `Makefile.in` in the `src` directory. I also made sure that my shell startup file included my CUDA executable and library paths, `/usr/local/cuda/bin` and `/usr/local/cuda/lib`.

I then ran `R CMD INSTALL` as above. I tested it by trying `gpuLm.fit()`, the `gputools` version of R's regular `lm.fit()`.

The package offers various linear algebra routines, such as matrix multiplication, solution of $Ax = b$ (and thus matrix inversion), and singular value decomposition, as well as some computation-intensive operations such as linear/generalized linear model estimation and hierarchical clustering.

Here for instance is how to find the square of a matrix `m`:

```
1 > m2 <- gpuMatMult(m,m)
```

The `gpuSolve()` function works like the R `solve()`. The call `gpuSolve(a,b)` will solve the linear system $ax = b$, for a square matrix `a` and vector `b`. If the second argument is missing, then a^{-1} will be returned.

1.8.3 rgpu 扩展包

为了安装 `rgpu`，我从 https://gforge.nbic.nl/frs/?group_id=38 下载源代码并解压缩。之后我修改了 `Makefile` 文件中的几行^⑩

```
1 LIBS = -L/usr/lib/nvidia -lcuda -lcudart -lcublas
2 CUDA_INC_PATH = /home/matloff/NVIDIA_GPU_Computing_SDK/C/common/inc
3 R_INC_PATH = /usr/include/R
```

第一行是为了让系统找到 `-lcuda`，这点和 `gputools` 一样。第二行是为了 NVIDIA SDK 中的 `cutil.h` 文件，上面的是我的安装路径。

For the third line, I made a file `z.c` consisting solely of the line

^⑩ 译者注：请根据自己机器上的相应路径进行修改

```
1 #include <R.h>
```

```
    and ran
```

```
1 R CMD SHLIB z.c
```

```
just to see whether the R include file was.
```

As of May 2010, the routines in **rgpu** are much less extensive than those of **gputools**. However, one very nice feature of **rgpu** is that one can compute matrix expressions without bringing intermediate results back from the device memory to the host memory, which would be a big slowdown. Here for instance is how to compute the square of the matrix **m**, plus itself:

```
1 > m2m <- evalgpu(m %*% m + m)
```

1.9 通过在 R 中调用 C 进行并行

Parallel R aims to be faster than ordinary R. But even if that aim is achieved, it's still R, and thus potentially slow.

One must always decide how much effort one is willing to devote to optimization. For the fastest code, we should not write in C, but rather in assembly language. Similarly, one must decide whether to stick purely to R, or go to the faster C. If parallel R gives you the speed you need in your application, fine; if not, though, you should consider writing part of your application in C, with the main part still written in R. You may find that placing the parallelism in the C portion of your code is good enough, while retaining the convenience of R for the rest of your code.

1.9.1 在 R 中调用 C

In C, two-dimensional arrays are stored in row-major order, in contrast to R's column-major order. For instance, if we have a 3x4 array, the element in the second row and second column is element number 5 of the array when viewed linearly, since there are three elements in the first column and this is the second element in the second column. Of course, keep in mind that C subscripts begin at 0, rather than at 1 as with R. In writing your C code to be interfaced to R, you must keep these issues in mind.

All the arguments passed from R to C are received by C as pointers. Note that the C function itself must return `void`. Values that we would ordinarily return must in the R/C context be communicated through the function's arguments, such as **result** in our example below.

1.9.2 Example: Extracting Subdiagonals of a Matrix

As an example, here is C code to extract subdiagonals from a square matrix.^① The code is in a file **sd.c**:

```
1 // arguments:
2 // m: a square matrix
3 // n: number of rows/columns of m
```

^①I wish to thank my former graduate assistant, Min-Yu Huang, who wrote an earlier version of this function.

```

4 // k: the subdiagonal index--0 for main diagonal, 1 for first
5 // subdiagonal, 2 for the second, etc.
6 // result: space for the requested subdiagonal, returned here
7
8 void subdiag(double *m, int *n, int *k, double *result)
9 {
10     int nval = *n, kval = *k;
11     int stride = nval + 1;
12     for (int i = 0, j = kval; i < nval-kval; ++i, j+= stride)
13         result[i] = m[j];
14 }

```

For convenience, you can compile this by running R in a terminal window, which will invoke GCC:

```

1 % R CMD SHLIB sd.c
2 gcc -std=gnu99 -I/usr/share/R/include -fpic -g -O2 -c sd.c -o sd.o
3 gcc -std=gnu99 -shared -o sd.so sd.o -L/usr/lib/R/lib -lR

```

Note that here R showed us exactly what it did in invoking GCC. This allows us to do some customization.

But note that this simply produced a dynamic library, **sd.o**, not an executable program. (On Windows this would presumably be a **.dll** file.) So, how is it executed? The answer is that it is loaded into R, using R's **dyn.load()** function. Here is an example:

```

1 > dyn.load("sd.so")
2 > m <- rbind(1:5, 6:10, 11:15, 16:20, 21:25)
3 > k <- 2
4 > .C("subdiag", as.double(m), as.integer(dim(m)[1]), as.integer(k),
5 result=double(dim(m)[1]-k))
6 [[1]]
7 [1] 1 6 11 16 21 2 7 12 17 22 3 8 13 18 23 4 9 14 19 24 5 10 15 20 25
8
9 [[2]]
10 [1] 5
11
12 [[3]]
13 [1] 2
14
15 $result
16 [1] 11 17 23

```

Note that we needed to allocate space for **result** in our call, in a variable we've named **result**. The value placed in there by our function is seen above to be correct.

1.9.3 在 R 中调用 OpenMPI C 代码

由于 OpenMP 可以由 C 使用，这就使得其可以从 R 中调用。（关于 OpenMP 的详细讨论请见第??章。）

在1.9节中类似，代码被编译并载入到 R 会话，尽管有一些额外的步骤用于在调用 GCC 时设置 `-fopenmp` 参数（你需要手动运行，而不是使用 `R CMD SHLIB`）。

1.9.4 在 R 中调用 CUDA 代码

The same principles apply here, but one does have to be careful with libraries and the like.

As before, we want to compile not to an executable file, but to a dynamic library file. Here's how, for the C file `mutlinksforr.cu` presented in the next section, the compile command is

```
1 pc41:~% nvcc -g -G -I/usr/local/cuda/include -Xcompiler
2   "-I/usr/include/R -fpic" -c mutlinksforr.cu -o mutlinks.o -arch=sm_11
3 pc41:~% nvcc -shared -Xlinker "-L/usr/lib/R/lib -lR"
4   -L/usr/local/cuda/lib mutlinks.o -o meanlinks.so
```

The product of this was `meanlinks.so`. I then tested it on R:

```
1 > dyn.load("meanlinks.so")
2 > m <- rbind(c(0,1,1,1),c(1,0,0,1),c(1,0,0,1),c(1,1,1,0))
3 > ma <- rbind(c(0,1,0),c(1,0,0),c(1,0,0))
4 > .C("meanout",as.integer(m),as.integer(4),mo=double(1))
5 [[1]]
6 [1] 0 1 1 1 1 0 0 1 1 0 0 1 1 1 1 0
7
8 [[2]]
9 [1] 4
10
11 $mo
12 [1] 1.333333
13
14 > .C("meanout",as.integer(ma),as.integer(3),mo=double(1))
15 [[1]]
16 [1] 0 1 1 1 0 0 0 0 0
17
18 [[2]]
19 [1] 3
20
21 $mo
22 [1] 0.3333333
```

1.9.5 示例: Mutual Outlinks

We again take as our example the mutual-outlinks example from Section ?? . Here is an R/CUDA version:

```

1 // CUDA example: finds mean number of mutual outlinks, among all pairs
2 // of Web sites in our set
3
4 #include <cuda.h>
5 #include <stdio.h>
6
7 // the following is needed to avoid variable name mangling
8 extern "C" void meanout(int *hm, int *nrc, double *meanmut);
9
10 // for a given thread number tn, calculates pair, the (i,j) to be
11 // processed by that thread; for nxn matrix
12 __device__ void findpair(int tn, int n, int *pair)
13 { int sum=0,oldsum=0,i;
14   for(i=0; ;i++) {
15     sum += n - i - 1;
16     if (tn <= sum-1) {
17       pair[0] = i;
18       pair[1] = tn - oldsum + i + 1;
19       return;
20     }
21     oldsum = sum;
22   }
23 }
24
25 // proclpair() processes one pair of Web sites, i.e. one pair of rows in
26 // the nxn adjacency matrix m; the number of mutual outlinks is added to
27 // tot
28 __global__ void proclpair(int *m, int *tot, int n)
29 {
30   // find (i,j) pair to assess for mutuality
31   int pair[2];
32   findpair(threadIdx.x,n,pair);
33   int sum=0;
34   // make sure to account for R being column-major order; R's i-th row
35   // is our i-th column here
36   int startrowa = pair[0],
37       startrowb = pair[1];
38   for (int k = 0; k < n; k++)

```

```

39     sum += m[startrowa + n*k] * m[startrowb + n*k];
40     atomicAdd(&tot,sum);
41 }
42
43 // meanout() is called from R
44 // hm points to the link matrix, nrc to the matrix size, meanmut to the output
45 void meanout(int *hm, int *nrc, double *meanmut)
46 {
47     int n = *nrc,msize=n*n*sizeof(int);
48     int *dm, // device matrix
49         htot, // host grand total
50         *dtot; // device grand total
51     cudaMalloc((void **)&dm,msize);
52     cudaMemcpy(dm,hm,msize,cudaMemcpyHostToDevice);
53     htot = 0;
54     cudaMalloc((void **)&dtot,sizeof(int));
55     cudaMemcpy(dtot,&htot,sizeof(int),cudaMemcpyHostToDevice);
56     dim3 dimGrid(1,1);
57     int npairs = n*(n-1)/2;
58     dim3 dimBlock(npairs,1,1);
59     proclpair<<<dimGrid,dimBlock>>>(dm,dtot,n);
60     cudaThreadSynchronize();
61     cudaMemcpy(&htot,dtot,sizeof(int),cudaMemcpyDeviceToHost);
62     *meanmut = htot/double(npairs);
63     cudaFree(dm);
64     cudaFree(dtot);
65 }

```

The code is hardly optimal. We should, for instance, have more than one thread per block.

1.10 调试 R 程序

R 内置的调试机制是首选，在还存在着其它选择。

1.10.1 文本编辑器

然而，如果你是一个 Vim 编辑器的粉丝，我开发了一个可以极大扩展 R 调试器的工具。请从 R 的 CRAN 上下载 `edtdbg`。Emacs 中也有类似的工具。

Vitalie Spinu 的 `ess-tracebug` 运行于 Emacs。它大体基于 `edtdbg`，但提供了更多的针对 Emacs 的特性。

1.10.2 IDE

我个人不是提倡使用 IDE，但的确有一些很优秀的 IDE。

REvolution Analytics, 一家提供 R 咨询和再开发版本 R 的公司, 他们提供了一个包含了很好的调试机制的 IDE。但它只可以在 Windows 上运行, 而且必须安装 Microsoft Visual Studio。

StatET, 一个基于 Eclipse 的跨平台 IDE 的开发者在 2011 年五月添加了调试工具。

RStudio, 另一个跨平台的 IDE 的开发, 从 2011 年夏天也开始计划添加调试器^②。

1.10.3 缺少命令行终端的问题

Parallel R packages such as **Rmpi**, **snow**, **foreach** and so on do not set up a terminal for each process, thus making it impossible to use R's debugger on the workers. What then can one do to debug apps for those packages? Let's consider **snow** for concreteness.

First, one should debug the underlying single-worker function, such as **mtl()** in our mutual outlinks example in Section 1.5.6. Here one would set up some artificial values of the arguments, and then use R's ordinary debugging facilities.

This may be sufficient. However, the bug may be in the arguments themselves, or in the way we set them up. Then things get more difficult. It's hard to even print out trace information, e.g. values of variables, since **print()** won't work in the worker processes. The **message()** function may work for some of these packages; if not, you may have to resort to using **cat()** to write to a file.

Rdsm allows full debugging, as there is a separate terminal window for each process.

1.10.4 调试 R 所调用的 C 代码

For parallel R that is implemented via R calls to C code, producing a dynamically-loaded library as in Section 1.9, debugging is a little more involved. First start R under GDB, then load the library to be debugged. At this point, R's interpreter will be looping, anticipating reading an R command from you. Break the loop by hitting ctrl-c, which will put you back into *GDB's* interpreter. Then set a breakpoint at the C function you want to debug, say **subdiag()** in our example above. Finally, tell GDB to continue, and it will then stop in your function! Here's how your session will look:

```
1 $ R -d gdb
2 GNU gdb 6.8-debian
3 ...
4 (gdb) run
5 Starting program: /usr/lib/R/bin/exec/R
6 ...
7 > dyn.load("sd.so")
```

1.11 本书中的其它 R 语言示例

见下列章节中的示例 (一些是非并行的):

??节、??节 (非并行) 和 ??节 (非并行)。

- 线性等式的并行 Jacobi 迭代, ??节。
- 1 维 FFT 的矩阵运算, ??节 (可以通过并行的矩阵相乘来并行化)。

^②译者注: RStudio 中的调试功能已添加

-
- 2 维 FFT 的并行计算，??节。
 - 图像平滑，??节。