

A Quick Guide for the pbdBASE Package

Drew Schmidt¹, Wei-Chen Chen², George Ostrouchov^{1,2},
Pragneshkumar Patel¹

¹Remote Data Analysis and Visualization Center
University of Tennessee,
Knoxville, TN, USA

²Computer Science and Mathematics Division,
Oak Ridge National Laboratory,
Oak Ridge, TN, USA

Contents

| | |
|---|------------|
| Acknowledgement | iii |
| Abstract | 1 |
| 1. Introduction | 1 |
| 1.1. Achievements | 1 |
| 1.2. Installation | 2 |
| 1.3. Package Examples | 2 |
| 2. Classes and Methods | 3 |
| 3. Information for Users | 4 |
| 3.1. The General Procedure for Using the System | 4 |
| 3.2. Printing in Parallel | 5 |
| 3.3. Process Grid Size | 5 |
| 3.4. Distributing Data | 6 |
| 3.5. Reading Data In Parallel | 6 |
| 4. Basic Example | 7 |
| 5. Using Information for Advanced Users | 9 |
| 5.1. Blocking Factor | 9 |
| 5.2. Different BLACS Contexts | 12 |
| 5.3. Exiting BLACS Contexts | 12 |
| 6. Advanced Example | 12 |

| | |
|--|-----------|
| 7. Information for Developers | 15 |
| 7.1. Class <code>ddmatrix</code> | 15 |
| 7.2. BLACS | 16 |
| References | 17 |

Acknowledgement

Ostrouchov, Schmidt, 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 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 thank our colleague, Ed D’Azevedo from the Computational Mathematics Group, Computer Science and Mathematics Division, Oak Ridge National Laboratory (ORNL), for his discussions and illuminating advice using ScaLAPACK and distributed matrix computation.

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

Abstract

With the size of data ever growing, the use of multiple processors in a single analysis becomes more and more a necessity. The Programming Big Data (pbd) project attempts to address the R language’s current shortcomings in parallel distributed computations. The **pbdBASE** package for R provides a distributed matrix datatype and low-level methods for this data type, including extraction via `[]`, as well as `NA` removal. Further, the package contains a set of BLACS, PBLAS, and ScaLAPACK wrappers. In addition to performance improvements through parallelism, use of this system with more than one processor allows the user to break R’s local memory barrier, namely the requirement that a vector be indexed by a 32-bit integer, by only storing subsets of the vector on each processor.

1. Introduction

The Programming with Big Data: BASE system, the R (R Core Team 2012) package **pbdBASE** (Schmidt *et al.* 2012a), is a (mostly) implicitly parallel foundational infrastructure to support higher level pbd packages, such as **pbdDMAT** (Schmidt *et al.* 2012b). Much of what it does is meant to live behind the scenes of packages further up the chain of the pbd ecosystem, and is largely targeted at developers. However, it does offer some essential functionality for all users.

In many ways, the **pbdBASE** package serves the pbd project in much the same way as R’s **base** package serves it. The principal goal of the **pbdBASE** package is to provide distributed classes (presently, a distributed dense matrix class), and many low-level functions for interacting with these classes. Many of these functions are wrappers of and for the distributed matrix algebra libraries BLACS, PBLAS, and ScaLAPACK. (Blackford *et al.* 1997) A set of S4 methods for R’s linear algebra functions using these wrappers is provided by a separate package, **pbdDMAT**.

Updates and bug releases for this and other **pbd** projects may, especially while in infancy, be much more frequent than CRAN releases. So for up to date packages, as well as evolving information about the **pbd** project, see the website “Programming with Big Data in R” at <http://r-pbd.org/>.

1.1. Achievements

The **pbdBASE** package, together with its sister packages in the **pbd** chain, offer the R user several new advancements. First, by making use of ScaLAPACK “under the hood”, we offer (near) ScaLAPACK speeds and scaling to many cores, but with R syntax. Second, by distributing the objects across processors, we are able to largely overcome R’s memory barrier.

At present¹, it is impossible to index native R objects with a 32-bit integer. Since a matrix in R is really just an array, this means that the largest square matrix it is possible to store in R is roughly a $46,000 \times 46,000$ matrix. This imposes two restrictions on the **pbd** system. First, the global dimension of any matrix used at this time with the **pbd** toolchain must have dimensions indexable by a 32-bit integer. Namely, no single dimension of the “full”, global

¹Though this is expected to change by summer 2013

matrix may have more than

$$(2^{32-1} - 1)^2 \approx 4.612 \times 10^{18}$$

because each dimension must be an integer, and in R terms, that means a 32-bit integer.

By comparison, the largest matrix which a single R process can hold has

$$\begin{aligned} 2^{32-1} - 1 &= 2,147,483,647 \\ &\approx 2 \times 10^9 \end{aligned} \tag{1}$$

numeric elements. However, we note that getting near the theoretical upper bound in (1) with the **pbd** system is effectively impossible, because each local R process will store at most roughly 10^9 elements. So even with 100,000 cores, you are still solidly within this boundary. Indeed, a user with N processors is able to store a square distributed matrix up to size

$$N \times (2^{32-1} - 1)$$

So at this time, a user would need 1024 cores to comfortably be able to analyze a terabyte of data, and over 100,000 cores to approach petabyte scale.

1.2. Installation

The **pbdBASE** package is available from the CRAN at <http://cran.r-project.org>, and can be installed via a simple

Installing pbdBASE

```
install.packages("pbdBASE")
```

This assumes only that you have MPI installed and properly configured on your system. If the user can successfully install the package's two principal dependencies, **pbdMPI** (Chen *et al.* 2012a) and **pbdSLAP** (Chen *et al.* 2012c) (each available from the CRAN), then the installation for **pbdBASE** should go smoothly. If you experience difficulty installing either these packages, you should see their documentation.

1.3. Package Examples

One can quickly get started with **pbdBASE** by learning from the following three examples:

```
### Under command mode, run the demo with 2 processors by
### (Use Rscript.exe for windows system)
mpiexec -np 2 Rscript -e "demo(example1, 'pbdBASE', ask=F, echo=F)"
mpiexec -np 2 Rscript -e "demo(example2, 'pbdBASE', ask=F, echo=F)"
mpiexec -np 2 Rscript -e "demo(example3, 'pbdBASE', ask=F, echo=F)"
```

1.4. Terminology

Before beginning, we will make frequent use of concepts from the Single Program/Multiple Data (SPMD) paradigm. If you are entirely unfamiliar with this approach to parallelism, or

if you are unfamiliar with the **pbdMPI** package, then you are strongly encouraged to read the vignette (Chen *et al.* 2012b) contained in the **pbdMPI** package, as well as examine and digest its many examples in order to better understand what follows.

A concise explanation of SPMD is that it is an approach to parallel, distributed programming in which one program is written, and each processor runs that same program, though that program locally will often be interacting with different data. This, in contrast to the manager/worker paradigm where one processor, the manager, is in charge of its workers, each of whom swear fealty to the manager. So in SPMD, each processor believes itself to be the manager, the one in charge. As a colleague, Dr. Russell Zaretzki put it, “it’s like academia.”

Throughout the remainder, we will be discussing distributed data objects such as matrices, and wish to do so with some standardized terminology. A matrix is of course a rectangular collection of numbers. A *distributed matrix* then is just a matrix which has been decomposed in some fashion so that each processor only owns a piece of the “whole” matrix. The “whole” matrix (which need not ever actually exist, except theoretically, at any time), rather than pieces of it distributed among the processors, will be referred to as a/the *global matrix*. Loosely speaking, the global matrix is what we are really thinking of when we deal with the distributed matrix.

In the SPMD paradigm, each processor, though only owning a piece of the whole (henceforth referred to as the *local matrix* or *submatrix*, relative to that processor), will call functions on that matrix exactly as one would with an ordinary, non-distributed matrix on a single processor. The difference for the user is minimal; all the “heavy lifting” which explicitly handles the distributed nature of the object is performed in the background.

Matrices, distributed or otherwise, have dimensions — that is, lengths of the number of rows and the number of columns in the rectangle. The global matrix has a *global dimension*, and this is a global value, i.e., this value does not vary from processor to processor. Every processor agrees as to the size of the “full” matrix, otherwise we would have anarchy. However, the local matrices, in practice, will differ from processor to processor, and so too should their *local dimensions*. A local dimension, as the name implies, is the dimension of the submatrix, relative to a particular processor.

2. Classes and Methods

Presently, package **pbdBASE** contains one new class, namely the class **ddmatrix** which stands for **d**istributed **d**ense **m**atrix. This S4 class serves as a container for a distributed matrix type, consisting of the members:

| | | |
|---|--------------|---|
| $\mathbf{ddmatrix} = \left\{ \begin{array}{ll} \end{array} \right.$ | Data | S4 slot containing the object’s submatrix, an R matrix |
| | dim | S4 slot containing the dimension of the global matrix, a numeric pair |
| | ldim | S4 slot containing the dimension of the local submatrix, a numeric pair |
| | bldim | S4 slot containing the ScaLAPACK blocking factor, a numeric pair |
| | CTXT | S4 slot containing the BLACS context, an numeric singleton |

with prototype

$$\text{new("ddmatrix")} = \begin{cases} \text{Data} & = \text{matrix}(0) \\ \text{dim} & = \text{c}(1,1) \\ \text{ldim} & = \text{c}(1,1) \\ \text{bldim} & = \text{c}(1,1) \\ \text{CTXT} & = 0 \end{cases}$$

We will discuss the last two items in more detail in the later sections, particularly Section 5 and Section 7.

In addition, the **pbdBASE** package contains new S4 methods for class `ddmatrix`, many of

| R S4 Method Overloading | New S4 Methods |
|--|--|
| [, [length(), dim(), nrow(), ncol() as.matrix(), as.vector() na.exclude() print() | submatrix(), submatrix<- ldim(), bldim(), ctxt() as.ddmatrix() |

Table 1: Class `ddmatrix` Methods in Package **pbdBASE**

which are listed in Table 1. Full information on these methods, as well as some new S3 methods and important non-method functions, is provided in the **pbdBASE** package documentation.

3. Information for Users

3.1. The General Procedure for Using the System

To use **pbdBASE**, and whence any package which relies on it, there are a few special considerations one must keep in mind which separate the system over using ordinary R. These are

1. In addition to the ordinary MPI communicator provided by **pbdMPI**, a special, rectangular MPI communicator is used. ([Dongarra and Whaley 1995](#))
2. Data is block cyclically distributed across the rectangular process grid. ([Blackford *et al.* 1997](#)) ([link](#))

The first item is fairly simple. Simply, you will start every analysis using the **pbdBASE** by making a call to the function `init.grid()`. The latter is slightly more complicated; for the moment, we will merely say that one achieves this with a parallel reader or data distribution function. See the package reference manual, as well as the later sections of this vignette for more details.

A generic skeleton of what a typical analysis using **pbdBASE** looks like involves 4 steps (which will be discussed in detail in the sections to follow):

1. Initialize the process grid. (`init.grid()`)
2. Read the data into the processes, store it as a distributed matrix.
3. Call R functions exactly as you would with ordinary matrices (when applicable).
4. Collect your results and finalize (`finalize()`).

One can generally use the above steps on existing R scripts with only a few minor modifications, and quickly parallelize his or her serial code. For most users, this will amount to simply adding in the appropriate calls to `init.grid()`, `finalize()`, and a parallel data reader or data distributor function (such as `as.ddmatrix()`). However, there are a few hitches with setting up a process grid and with distributing/parallel reading. We discuss these issues in the following two sections.

3.2. Printing in Parallel

Printing in parallel can take some getting used to, especially in SPMD style programs. If you simply issue the order `print(x)`, then every process will print, but it is often a cacophony of undecipherable writing out to the terminal, with each process trampling the others.

To this end, you should learn to make extensive use of the **pbdMPI** function `comm.print()`. For instance, if you have two processors, each with an object `x`, but with disagreement to what `x` actually is — say for example, one thinks `x` is 1, and the other thinks it is 2. Then calling `comm.print(x)` will print two pieces of information: the processor printing the value (here process 0) and the value itself, 1. By default, the only processor to print is that with MPI communicator address 0. You can have all processor ranks print using the optional argument `all.rank=TRUE`. This will print all values stored for the requested object, but will do so “one at a time”. You can also disable the printing of which processor is doing the printing via the optional argument `quiet=TRUE`. See the official **pbdMPI** documentation for details.

Additionally, there is a method for the `print()` function when applied to a distributed class such as `ddmatrix`. This will print some brief information about the matrix from processor 0. This has the optional argument `all=`, which can be used to print the entire matrix, one line at a time.

This function should not be called from `comm.print()`. Actually, one of the easiest ways to get yourself into trouble and hang up all the processors is to call a function which requires communication between processors from inside something like `comm.print()`.

3.3. Process Grid Size

Recall that we will very frequently visualize the processors as being in a 2-dimensional processor grid. This grid is initialized via the function `init.grid()`, which accepts the optional arguments `nprows=` and `npcols=`. If these are left blank, then a reasonable choice will be made based on the number of available processors. Here “reasonable” means “as close to square as possible.” The inspired reader can find more detail within the ScaLAPACK User’s Guide (Blackford *et al.* 1997) as to why this is a good choice. ([link](#))

To reiterate, in most cases, taking `nprow` and `npcol` as close to each other as possible is “sufficiently good”. Leaving the `nprow=` and `npcol=` options blank will make this choice for you. For example, we note that the user should be aware that providing 37 cores to **pbdDMAT** may not perform as well as providing 36 cores in the form of a 6×6 process grid. There is a strong connection between the process grid and the block-cyclic distribution, which we will discuss further in the following section.

Additionally, we note that the `init.grid()` function accepts an additional argument `ICTXT=`, which we will discuss in Section 7. This argument is also discussed in detail in the package reference manual.

3.4. Distributing Data

When distributing data, you must use a *blocking factor*. This is a pair of numbers (a, b) , and unless you think you have a great reason to do otherwise, you should have $a = b$. If you have no intuition, the just make them equal. The scale of these numbers should generally correspond to the size of your process grid and the scale of the data. The choice of blocking factor can seriously impact performance, because it is intimately tied to the data distribution. We will spare the details for the moment, and merely say that the blocking factor constitutes a tradeoff. Smaller values, down to $(1, 1)$, mean that there will be more parallelism in many of the matrix algebra routines but also increase communication costs. On the other hand, larger values, which could be larger than the dimensions of any one of your matrices, will limit communication between the processors, but naturally also limit the parallelism.

A good choice of blocking factor can be difficult, and in the author’s opinion, requires some intuition gained through experimentation. A discussion on this topic can be found in the ScaLAPACK User’s Guide. (Blackford *et al.* 1997) ([link](#)) However, we note that **pbdBASE** defaults to a 4×4 blocking factor, which is probably an acceptable choice if you do not know what else to do. For monstrously huge matrices, it could be a bit small, and so scaling it up by powers of 2 (8×8 , 16×16 , \dots , 256×256) may be warranted. The blocking does not have to be by power of 2, but this is a convenient way to do business. The performance-hungry user is encouraged to experiment with various blocking factors across various processor grids with various matrix sizes to develop intuition.

3.5. Reading Data In Parallel

To really get the most out of this system, you need to read the data into R in a parallel distributed fashion. This generally necessitates the use of a parallel file system, such as Lustre. These are the kinds of resources that one generally does not have on his/her laptop, unfortunately. It is possible to read all of the data in on one core and have that core distribute the data to all the other processes, which is what one should do in the absence of a parallel file system. However, these added communication costs could overtake the gains provided by distributed computation, depending on the task. Worse, the user is again trapped in the world of 32-bit integer indexing, meaning the size of problem that it is possible to solve shrinks.

As a general rule, if you are on a smaller system with limited resources, you do what you must. If you are on a larger system with luxurious resources, you really should know better, and act accordingly. Failure to do so will significantly negatively impact performance with

this system.

Some functions which are useful in this regard are `as.ddmatrix()`, `distribute()`, and `redistribute()`. See the reference manual for details.

4. Basic Example

Now that you've come to this section directly from the abstract, skipping the other sections, let's take a look at an example. Since the **pbdbase** package is mostly about providing functionality to packages farther up the pbd chain, it won't be particularly exciting. But you have to crawl before you can drag race.

First we will generate some data on process 0, or (0,0) using the grid notation. We will do this using a simple `if()` together with the **pbdMPI** function `comm.rank()`, which returns the MPI communicator number for the calling process. Any process not initially storing any data should store `NULL` for the object. This probably isn't the way you will want to run your production code, especially if you are randomly generating data; in that case, it would be much more efficient to just generate what is needed on each processor. Although doing this requires that you have access to good seeds for parallel random number generation; for more information, see the documentation on setting seeds via `comm.set.seed()` in the **pbdMPI** package, or for more control, see the **rlecuyer** or **rsprng** packages.

There is merit, however, in operating in this way. This is somewhat like the process necessary for reading in data onto a subset of processors (just 1 if you do not have access to a parallel file system) and then distributing that out to the larger grid, so it is a useful skill. For more information about this procedure, see Section 6.

Generating Test Data

```
library(pbdbase, quiet=TRUE)

init.grid()

nrows <- ncols <- 10 # generate a 10x10 matrix
BL <- 2 # the blocking factor for distribution

# Generate data on process 0, then distribute to the others
{
  if (comm.rank()==0)
    x <- matrix(rnorm(n=nrows*ncols, mean=10, sd=100),
               nrow=nrows, ncol=ncols)
  else
    x <- NULL
}

dx <- as.ddmatrix(x=x, bldim=BL)

# continued in the next block of code ...
```

To convince ourselves that the data is distributed, we can inspect the new object in several ways:

Printing the Object

```
print(dx)

comm.print(submatrix(dx))

comm.print(dx)

# continued in the next block of code ...
```

Here, `print()` is a special method that shows you the slots of your distributed matrix. The `submatrix()` function will show the local submatrix (syntactic sugar for printing `dx@Data`). Use of `pbdMPI`'s `comm.print()` ensures that only process 0 will print the result. Finally, just using R's `print` method on the object in `comm.print(dx)` will produce an uglier version of `print(dx)` and `comm.print(submatrix(dx))`.

We can also do insertions and extractions:

Insertion and Extraction

```
dx[1,1] <- NA # insertion indices are global
comm.print(submatrix(dx)[1,1], all.rank=T) # see?

comm.print(dim(dx))
dx <- dx[, -2]
comm.print(dim(dx))

nona <- na.exclude(dx)

# continued in the next block of code ...
```

Finally, we can convert the distributed matrix back into an ordinary R matrix on processor 0. You probably will not need to do this very often in production code, because in practice, you could be dealing with matrices with so many elements that they will not fit into a single R process. For testing however, this process can be very useful. It could also conceivably have utility for dealing with $n \times 1$ matrices.

Insertion and Extraction

```
# convert back
nona <- as.matrix(nona, proc.dest=0)

# compare our results with R --- notice the syntax is
# essentially identical
if (comm.rank()==0){
  x[1,1] <- NA
```

```

x <- x[, -2]
r_nona <- na.exclude(x)

all.equal(r_nona, nona)
}

finalize()

```

In the above script, there is one addition over the previous pieces. Namely, we include several calls to `comm.cat()`. All this does is demand line breaks (via the regular expression `\n`) for more human-readable printing.

The script file is available in the **pbdBASE** directory, under `inst/examples/base_eg.R`, and you can run this script from the command line with the following command:

```

# replace the 4 below with your number of processors
mpirun -np 4 Rscript pbdbase_example.R

```

If you want to ramp up the size of the problem and the number of cores, you may want to change the `nrows`, `ncols`, and BL definitions (these are good to experiment with regardless).

5. Information for Advanced Users

In this section, we will discuss the block-cyclic distribution of data across a 2-dimensional processor grid in lengthy detail. It probably goes without saying that before beginning this section, the reader should be familiar with all sections prior.

5.1. Blocking Factor

The motivation for the development of **pbdBASE** was to be able to use ScaLAPACK routines for distributed linear algebra. The S4 methods which use these routines are in the package **pbdDMAT**, but all the same, using block-cyclically distributed data is absolutely essential for **pbdBASE**, since it is essential for ScaLAPACK.

Most of the difficulty in using this type of distributed data has been abstracted away for the user of **pbdBASE**. However, you may still find it useful to experiment with a block cyclic data distributor provided by ([The ScaLAPACK Team 2010](#)) to gain some intuition and understanding of how the package is powered underneath, and why a simple `read.table()` call will not suffice.

The idea is fairly simple, even if the execution can sometimes be cumbersome. We want to try to evenly balance the data distribution for “large” matrices across the process grid, but in a way that is congruent with the natural blocking of matrices when performing LAPACK operations on them. As the name implies, we imagine taking a large global matrix and chopping it into blocks, and assigning those blocks to the processors in the grid. The way this is done is far from arbitrary, however.

For simplicity, let us explicitly assume for the moment that we are going to distribute a 9×9

matrix across a 2×3 process grid, using a 2×2 blocking factor. The above link may be extremely useful to your understanding the following. For those who flatly refuse to click the link, we can visualize the process as follows. We can imagine our global matrix (even if we never actually have a globally stored matrix, this is the way we would imagine it in our heads) as looking like:

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

and our process grid looks like:

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

with the usual MPI processor rank on the left, and the corresponding BLACS processor grid position on the right.

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.

Two visual representations of this block cyclic data decomposition can be seen in Figure 1 and Figure 2 . As you can see “all animals are equal, but some animals are more equal than others.” Meaning that we don’t inherently penalize any processor or go out of our way to imbalance the data load, but that it is possible to do so with poor choice of blocking factor, and that generally process $(0, 0)$ in the BLACS grid will receive the most data. Additionally, notice that matrices that are closer to being square will distribute better than will vectors. In this case, a vector (here, a 9×1 matrix) would only distribute across the processes in the first process column, namely $(0, 0)$ and $(1, 0)$. This isn’t necessarily the great imbalance problem you might instinctively think it is. These matrices are generally fairly small, so it’s not really that big of a deal that they do not distribute well. On the other hand, things close to square are really quite large, and so we should go out of our way to make sure that they distribute “well”, whatever that really means.

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

Figure 1: Block Cyclic Data Decomposition Example

$$\begin{array}{ccc} \begin{bmatrix} x_{11} & x_{12} & x_{17} & x_{18} \\ x_{21} & x_{22} & x_{27} & x_{28} \\ \hline x_{51} & x_{52} & x_{57} & x_{58} \\ x_{61} & x_{62} & x_{67} & x_{68} \\ \hline x_{91} & x_{92} & x_{97} & x_{98} \end{bmatrix}_{5 \times 4} & \begin{bmatrix} x_{13} & x_{14} & x_{19} \\ x_{23} & x_{24} & x_{29} \\ \hline x_{53} & x_{54} & x_{59} \\ x_{63} & x_{64} & x_{69} \\ \hline x_{93} & x_{94} & x_{99} \end{bmatrix}_{5 \times 3} & \begin{bmatrix} x_{15} & x_{16} \\ x_{25} & x_{26} \\ \hline x_{55} & x_{56} \\ x_{65} & x_{66} \\ \hline x_{95} & x_{96} \end{bmatrix}_{5 \times 2} \\ \begin{bmatrix} x_{31} & x_{32} & x_{37} & x_{38} \\ x_{41} & x_{42} & x_{47} & x_{48} \\ \hline x_{71} & x_{72} & x_{77} & x_{78} \\ x_{81} & x_{82} & x_{87} & x_{88} \end{bmatrix}_{4 \times 4} & \begin{bmatrix} x_{33} & x_{34} & x_{39} \\ x_{43} & x_{44} & x_{49} \\ \hline x_{73} & x_{74} & x_{79} \\ x_{83} & x_{84} & x_{89} \end{bmatrix}_{4 \times 3} & \begin{bmatrix} x_{35} & x_{36} \\ x_{45} & x_{46} \\ \hline x_{75} & x_{76} \\ x_{85} & x_{86} \end{bmatrix}_{4 \times 2} \end{array}$$

Figure 2: Local Processor Storage View of the Block Cyclic Data Decomposition

Now, the user should never need to develop an algorithm to perform this kind of block-cyclic decomposition of the data; it may be necessary, especially for a developer who wishes to use this system, to have to deal with this data distribution business in a very up-front way, but many helper functions are provided by the package to that end. We bring up this issue to illustrate a critical point: the choice of blocking can make a large difference in performance, and generally should be tailored to the process grid.

If the blocking factor (in the above, 2×2) is too “big” (relative to the process grid), then the data distribution will be very uneven. This has the net effect of reducing communication times between processes, but also limiting the amount of parallelism possible. On the other hand, if the blocking factor is too “small”, then the data distribution is “too fair”, gaining much parallelism but massively inflating communication costs. At the extreme ends of this are, for the former, using a blocking factor that would encompass the entirety of the matrix (so that the data is distributed to only one process — no communication, no parallelism), and for the latter, using a blocking factor of $(1, 1)$, so that there is effectively nothing *but* communication and parallelism.

Very loosely speaking, parallelism is good and communication is bad, but in practice, we can’t really have one without the other.

5.2. Different BLACS Contexts

A very advanced user, perhaps someone already familiar with ScaLAPACK, may wish to be able to use a variety of BLACS contexts. Herein, we shall discuss in depth the BLACS system and **pbdBASE**'s relationship with it.

Suppose you have np processors. When a call to `init.grid()` is first made, contexts 0, 1, and 2 are created. These are, respectively, the $P \times Q$ processor grid, where $P = nprows$ and $Q = npcols$ (with an optimal choice made for each of P and Q if the `nprows=` and `npcols=` arguments are missing), 1 is the $1 \times np$ process grid, and 2 is the $np \times 1$ process grid. For this reason, **contexts 0, 1, and 2 are reserved**. Additional contexts can be created via `init.grid()` by passing integer values of 3 or greater to the `ICTXT=` argument. The function `minctxt()` returns the smallest integer which is not being used as a BLACS context.

Throughout **pbdBASE** (and **pbdDMAT**), if a function does not accept an argument for BLACS context, you can rest assured that it does not alter context. However, there are some functions which will implicitly change context, such as `na.exclude()`, and some explicitly such as `redistribute()`. For this particular function, the data is redistributed across context 1, because in this distribution, dropping rows (namely those containing NA's) does not destroy block-cyclicalilty. For this reason, an optional context argument is provided, and the default will simply convert the matrix back to its incoming context (from the slot `@CTXT`).

5.3. Exiting BLACS Contexts

You can create BLACS contexts, and you can also free them. To free a particular context, you would use the `gridexit()` function. You can not free contexts 0, 1, or 2, as these are protected values that much of the internals rely on for the vast majority of users. You can free any other context number, assuming that you created one.

You can also shut down all BLACS contexts at once without shutting down all MPI communicators using the `blacsexit()` function. This is automatically called, as necessary, when you call `finalize()`, because a failure to shut down all BLACS communicators can result in memory leaks.

The names of these functions are the same as the names of the corresponding BLACS routines that they rely on.

6. Advanced Example

In this section, we will look at an example of how to create and utilize additional BLACS contexts, and make a pretend parallel reader.

Let's take a look at an example using the ideas discussed in Section 5. We will again be using 4 processors in the form of a 2×2 grid. However, this time, we are going to randomly generate a matrix on 2 processors and then distribute that data onto the full 2×2 grid. We will achieve this by creating a new BLACS context with 2 processor rows and 1 processor column. Processors not in this grid will have the BLACS communicator grid location of $(-1, -1)$.

For clarity, because this can be a confusing and deeply frustrating concept at first. Every processor is part of the BLACS context; *all of them*. The context is just a collection of information that makes it possible to perform local operations with global communications. Think of it like an R list; it's just a place to put things. Each processor stores its own copy of this "list", with some pieces common to all processors and some differing from process to process. The objects within that list are the BLACS context number (just an identifying name more so than actual information), the number of global processor rows/columns, and that processor's position in that grid. Whenever a function relying on BLACS (like all PBLAS and ScaLAPACK functions) is called, a context is required. The local processor's grid location value of $(-1, -1)$ when that processor is not *really* part of the grid is the BLACS equivalent of NULL. It's just a placeholder that lets the local processor know that it isn't part of the party. But that processor is still aware of the BLACS context, which is a global number, and the total number of processor rows/columns (also global). So it is the location in the processor grid that is $(-1, -1)$, not the context. The context is just a "name", which happens to be a global integer common to all processors.

The choice of 1 processor column is not random. Doing so demands that contiguous rows are stored on the processes, and the cyclic distribution is occurring between rows. This is a very convenient way to do business if you must read in data from files (rather than randomly generate it in memory). Here, we can think of each processor "reading in" the parts of the (imaginary) file, and then distributing that data out to BLACS grid 0 for analytics. Other, "non-vector" BLACS grids are certainly possible as well.

A word of caution; the seeds here on the different processors are not guaranteed to be independent. This is just for the sake of demonstration. Refer back to the comments in Section 4.

Generating in Parallel

```
library(pbdBASE, quiet=TRUE)

init.grid(nprow=2, npcol=2)

# find the minimum value possible for a new BLACS context
# the value should be 3
newctxt <- minctxt()
comm.print(newctxt)

# create new grid
init.grid(nprow=2, npcol=1, ICTXT=newctxt)

# store new grid information
grid <- blacs(ICTXT=newctxt)

# "read in" the data and distribute
if (grid$MYROW == -1 || grid$MYCOL == -1){
  x <- matrix(0)
} else {
  x <- matrix(rnorm(50), ncol=10)
```



```

}

dx <- new("ddmatrix",
          Data=x,
          dim=c(10,10), ldim=dim(x), bldim=c(5,10),
          CTXT=newctxt)

dx <- redistribute(dx, bldim=2, ICTXT=0)

print(dx)

# close grid
gridexit(newctxt)

blacsexit()

comm.print("MPI still works")

finalize()

```

A few comments. First, we do not technically need to call `gridexit()`, since the following `blacsexit()` call will shut down *all* BLACS grids. However, inbetween those two calls, we could make more calls to routines using BLACS grids 0, 1, and 2. Second, if we explicitly call `blacsexit()` as we do here, then all BLACS grids are shut down without shutting down the MPI communicator. However, `finalize()` will do this for us if we forget, because a failure to do so can cause memory leaks. So you do not have to explicitly call `blacsexit()` unless you are done with BLACS, have yet more MPI work to do, and want to free up the resources.

Additionally, notice that here we get friendly with the `new` constructor. This call is part of R's S4 methods, and instantiates a — as the name implies — new object of specified class with specified slots. Notice that the blocking dimension is the dimension of the local matrix. This is so because we are imagining reading in large, contiguous blocks for each processor (here with just 1 cycle). This is fairly ad hoc, but is useful for demonstration purposes. A more advanced example, which generates only what is needed on each processor by making use of the function `pbdBASE::numroc()` can be found in the vignette for **pbdDMAT**.

Finally, you may be wondering why we would even bother with this approach with contexts rather than simply explicitly choosing a subset of processors from context 0. We could do this as well, but this isn't quite as simple as you might think, especially with the tools already built (meaning you may have to work much, much harder for this). You are encouraged to simply construct a new BLACS context as in the example, because for this very low-level data wrangling, it can make your life much simpler.

Finally, this script is available in the **pbdBASE** directory, under `inst/examples/base_eg.R` and is meant to be run with 4 processors.

7. Information for Developers

In this section, we will discuss some issues that are not really important for anyone except those who need to develop new methods which rely on **pbdBASE**. For the remainder, we will make mention and use of R's S4 method of object oriented programming, and we will assume that the reader has at least a working familiarity with S4. There are several fine introductions to S4 methods available, including those from (Chambers 2006) and (Genolini 2008). It probably goes without saying that before beginning this section, the reader should be familiar with all sections prior.

7.1. Class `ddmatrix`

Every local submatrix must be a matrix containing at least one value, even if it technically should not have anything. You will unleash unspeakable misery on yourself if you use `matrix(nrow=0, ncol=0)` for the submatrix, in particular when passing down to ScaLAPACK routines (the motivating example for the existence of **pbdBASE**). Instead, you should use the routine `ownany()`, a simple wrapper on `numroc()`, to determine if the process actually owns any data from the global matrix. If you consider this an extravagant waste of memory, then I have some very interesting things to tell you about R.

Likewise, the slot `@ldim` is the dimension of the local storage, and not, necessarily, the “actual” local dimension (which could effectively be `NULL`). So for example, if you imagine all data for a distributed matrix living on process 0, then all the other processes should store `matrix(0)` in the `@Data` slot, and `c(1,1)` in the `@ldim` slot.

Let's take a look at an example. Suppose we want to develop a way of taking logs of the entries of our distributed matrix in a way that is a natural extension to that of R (this is actually already done in **pbdDMAT**, but is a good illustration). A quick call to `isGeneric(f="log")` shows that the function `log()` in R is already S4 generic, so we can easily enough define a new method for it.

For some problems, some processors will own `matrix(0)` under the `@Data` slot when really, technically, they store nothing. We don't want to take the log of these zeros, since there is no real point. Worse, failure to exclude these values across several methods can accumulate in incorrect answers (think of taking the log on all submatrices and then running a `sum()` method on a distributed matrix). Using `ownany()`, this is trivial.

Generating in Parallel

```
mylog <- function(x, base=exp(1))
{
  if (ownany(dim=x@dim, bldim=x@bldim, CTEXT=x@CTEXT))
    x@Data <- log(x=x@Data, base=base)
  return(x)
}
```

Now we just need to set the method:

Generating in Parallel

```
setMethod("log", signature(x="ddmatrix"),  
  mylog  
)
```

This is more or less how things are done in **pbdDMAT**.

7.2. BLACS

The most critical piece of information to impart is that no matter what, “block-cyclicity” can not be destroyed. Most of the routines for distributed matrices (and almost all of the really complicated ones) assume this structure. This is why 3 BLACS contexts are created by default when initializing the process grid; namely, context 0 is the optimal (unless otherwise specified) rectangular process grid, context 1 is the $PQ \times 1$ process grid (assuming P process rows and Q process columns), and context 2 is the $1 \times PQ$ process grid. As defined above, context 0 is optimal for many linear algebra routines, and not really any worse (aside from being cumbersome) for most other computations. The other contexts are important because they are more natural for adding/removing columns and rows (respectively).

Redistributing the data from a $P \times Q$ process grid to a $1 \times PQ$ process grid can be very useful, despite the communication overhead involved. When a matrix is distributed across this context, while it is not trivial to “drop” rows (certainly not as easy as it is in context 1), doing so, whichever rows we do indeed drop, results in a block-cyclically distributed matrix. In general, the same cannot be said for other contexts. The **ddmatrix** methods `[` and `na.exclude()` rely on such a redistribution (two redistributions in the case of the former).

References

- Blackford LS, Choi J, Cleary A, D'Azevedo E, Demmel J, Dhillon I, Dongarra J, Hammarling S, Henry G, Petitet A, Stanley K, Walker D, Whaley RC (1997). *ScaLAPACK Users' Guide*. Society for Industrial and Applied Mathematics, Philadelphia, PA. ISBN 0-89871-397-8 (paperback). URL http://netlib.org/scalapack/slug/scalapack_slug.html/.
- Chambers J (2006). "How S4 Methods Work." *Technical report*, The R-Project for Statistical Computing. URL <http://developer.r-project.org/howMethodsWork.pdf>.
- Chen WC, Ostrouchov G, Schmidt D, Patel P, Yu H (2012a). "pbdMPI: Programming with Big Data – Interface to MPI." R Package, URL <http://cran.r-project.org/package=pbdMPI>.
- Chen WC, Ostrouchov G, Schmidt D, Patel P, Yu H (2012b). "A Quick Guide for the pbdMPI package." R Vignette, URL <http://cran.r-project.org/package=pbdMPI>.
- Chen WC, Schmidt D, Ostrouchov G, Patel P (2012c). "pbdSLAP: Programming with Big Data – Scalable Linear Algebra Packages." R Package, URL <http://cran.r-project.org/package=pbdSLAP>.
- Dongarra J, Whaley RC (1995). "A User's Guide to the BLACS." *Technical report*, University of Tennessee. UT-CS-95-281, URL <http://www.netlib.org/lapack/lawnspdf/lawn94.pdf>.
- Genolini C (2008). "A (Not So) Short Introduction to S4." *Technical report*, The R-Project for Statistical Computing. URL <http://cran.r-project.org/doc/contrib/Genolini-S4tutorialV0-5en.pdf>.
- R Core Team (2012). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria. ISBN 3-900051-07-0, URL <http://www.r-project.org/>.
- Schmidt D, Chen WC, Ostrouchov G, Patel P (2012a). "pbdBASE: Programming with Big Data – Core pbd Classes and Methods." R Package, URL <http://cran.r-project.org/package=pbdBASE>.
- Schmidt D, Chen WC, Ostrouchov G, Patel P (2012b). "pbdDMAT: Programming with Big Data – Distributed Matrix Algebra Computation." R Package, URL <http://cran.r-project.org/package=pbdDMAT>.
- The ScaLAPACK Team (2010). "Block Cyclic Data Distribution." URL <http://acts.nersc.gov/scalapack/hands-on/datadist.html/>.