DMAT
○○○○○○○○○
○○○○○○
○○○○○○○○

pbdDMAT eg's
○○○○
○○○

# Programming with Big Data in R

Drew Schmidt and George Ostrouchov

August 8, 2013

**DMAT**
○○○○○○○○○
○○○○○○
○○○○○○○○

**pbdDMAT eg's**
○○○○
○○○

## About This Presentation

### Downloads

This presentation and supplemental materials are available at:

> http://r-pbd.org/tutorial

Sample R scripts and pbs job scripts available on Chester:
> /lustre/scratch/sw/r/3.0.1.new/chester/gnu4.7.3/
> EXAMPLES/scripts.tar.gz

**DMAT**
○○○○○○○○○
○○○○○○
○○○○○○○○

**pbdDMAT eg's**
○○○○
○○○

# About This Presentation

## Speaking Serial R with a Parallel Accent

The content of this presentation is based in part on the
**pbdDEMO** vignette *Speaking Serial R with a Parallel Accent*

http://goo.gl/HZkRt

It contains more examples, and sometimes added detail.

DMAT
○○○○○○○○○
○○○○○
○○○○○○○○

pbdDMAT eg's
○○○○
○○○

# About This Presentation

## Installation Instructions

Installation instructions for setting up a pbdR environment are available:

$$\texttt{http://r-pbd.org/install.html}$$

This includes instructions for installing R, MPI, and pbdR.

DMAT
○○○○○○○○
○○○○○○
○○○○○○○○

pbdDMAT eg's
○○○○
○○○

# Contents

DMAT
○○○○○○○○○
○○○○○○
○○○○○○○○

pbdDMAT eg's
○○○○
○○○

# Contents

1. Introduction to pbdDMAT and the DMAT Structure
   - Introduction to Distributed Matrices
   - DMAT Distributions
   - pbdDMAT

**DMAT**
●○○○○○○○
○○○○○○
○○○○○○○○

**pbdDMAT eg's**
○○○○
○○○

**Introduction to Distributed Matrices**

## Distributed Matrices

Most problems in data science are matrix algebra problems, so:

$$\text{Distributed matrices} \implies \text{Handle Bigger data}$$

**DMAT**
○●○○○○○○
○○○○○○
○○○○○○○○○

**pbdDMAT eg's**
○○○○
○○○

Introduction to Distributed Matrices

## Distributed Matrices

High level OOP allows *native* serial R syntax:

```
1  x <- x[-1, 2:5]
2  x <- log(abs(x) + 1)
3  xtx <- t(x) %*% x
4  ans <- svd(solve(xtx))
```

However...

**DMAT**
○○●○○○○○
○○○○○○
○○○○○○○○

**pbdDMAT eg's**
○○○○
○○○

**Introduction to Distributed Matrices**

### Distributed Matrices

DMAT:

- **D**istributed **MAT**rix data structure.
- No single processor should hold all of the data.
- Block-cyclic matrix distributed across a 2-dimensional grid of processors.
- Very robust, but confusing data structure.

DMAT
○○○●○○○○
○○○○○○
○○○○○○○○

pbdDMAT eg's
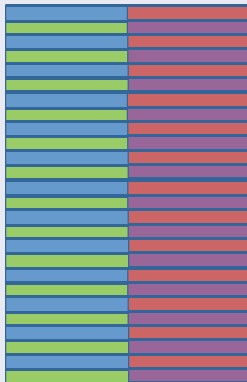○○○○
○○○

Introduction to Distributed Matrices

## Distributed Matrices



(a) Block     (b) Cyclic     (c) Block-Cyclic

Figure: Matrix Distribution Schemes

DMAT
○○○○●○○○
○○○○○○
○○○○○○○○

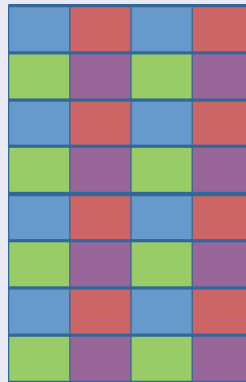pbdDMAT eg's
○○○○
○○○

Introduction to Distributed Matrices

## Distributed Matrices



(a) 2d Block    (b) 2d Cyclic    (c) 2d Block-Cyclic

Figure: Matrix Distribution Schemes Onto a 2-Dimensional Grid

**DMAT**
○○○○○●○○
○○○○○○
○○○○○○○○

**pbdDMAT eg's**
○○○○
○○○

Introduction to Distributed Matrices

## Processor Grid Shapes

$$\begin{bmatrix} 0 \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{bmatrix}^{T} \qquad \begin{bmatrix} 0 & 1 & 2 \\ 3 & 4 & 5 \end{bmatrix} \qquad \begin{bmatrix} 0 & 1 \\ 2 & 3 \\ 4 & 5 \end{bmatrix} \qquad \begin{bmatrix} 0 \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{bmatrix}$$

(a) $1 \times 6$      (b) $2 \times 3$      (c) $3 \times 2$      (d) $6 \times 1$

Table: Processor Grid Shapes with 6 Processors

**DMAT**
○○○○○○○●○
○○○○○○
○○○○○○○○

**pbdDMAT eg's**
○○○○
○○○

**Introduction to Distributed Matrices**

## Distributed Matrices

The data structure is a special R class (in the OOP sense) called ddmatrix. It is the "under the rug" storage for a block-cyclic matrix distributed onto a 2-dimensional processor grid.

$$\text{ddmatrix} = \begin{cases} \textbf{Data} & \text{S4 local submatrix, an R matrix} \\ \textbf{dim} & \text{S4 dimension of the global matrix, a numeric pair} \\ \textbf{ldim} & \text{S4 dimension of the local submatrix, a numeric pair} \\ \textbf{bldim} & \text{S4 ScaLAPACK blocking factor, a numeric pair} \\ \textbf{CTXT} & \text{S4 BLACS context, an numeric singleton} \end{cases}$$

with prototype

$$\text{new("ddmatrix")} = \begin{cases} \textbf{Data} & = \texttt{matrix(0.0)} \\ \textbf{dim} & = \texttt{c(1,1)} \\ \textbf{ldim} & = \texttt{c(1,1)} \\ \textbf{bldim} & = \texttt{c(1,1)} \\ \textbf{CTXT} & = \texttt{0.0} \end{cases}$$

**DMAT**
○○○○○○○●
○○○○○○
○○○○○○○○

**pbdDMAT eg's**
○○○○
○○○

Introduction to Distributed Matrices

## Distributed Matrices: The Data Structure

Example: an $9 \times 9$ matrix is distributed with a "block-cycling" factor of $2 \times 2$ on a $2 \times 2$ processor grid:



$$= \begin{cases} \textbf{Data} & = \texttt{matrix(...)} \\ \textbf{dim} & = \texttt{c(9, 9)} \\ \textbf{ldim} & = \texttt{c(...)} \\ \textbf{bldim} & = \texttt{c(2, 2)} \\ \textbf{CTXT} & = 0 \end{cases}$$

See http://acts.nersc.gov/scalapack/hands-on/datadist.html

**DMAT**
⊙⊙⊙⊙⊙⊙⊙⊙
●⊙⊙⊙⊙⊙
⊙⊙⊙⊙⊙⊙⊙⊙
**DMAT Distributions**

**pbdDMAT eg's**
○○○○
○○○

## Understanding Dmat: Global Matrix

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

DMAT
○○○○○○○○
○●○○○○
○○○○○○○○○

pbdDMAT eg's
○○○○
○○○

DMAT Distributions

## DMAT: 1-dimensional Row Block

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

$$\text{Processor grid} = \begin{vmatrix} 0 \\ 1 \\ 2 \\ 3 \end{vmatrix} = \begin{vmatrix} (0,0) \\ (0,1) \\ (1,0) \\ (1,1) \end{vmatrix}$$

**DMAT**
○○○○○○○○○
○○●○○○
○○○○○○○○○
**DMAT Distributions**

**pbdDMAT eg's**
○○○○
○○○

## DMAT: 2-dimensional Row Block

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

$$\text{Processor grid} = \begin{vmatrix} 0 & 1 \\ 2 & 3 \end{vmatrix} = \begin{vmatrix} (0,0) & (0,1) \\ (1,0) & (1,1) \end{vmatrix}$$

DMAT
○○○○○○○○○
○○○●○○
○○○○○○○○○○

pbdDMAT eg's
○○○○
○○○

DMAT Distributions

## DMAT: 1-dimensional Row Cyclic

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

$$\text{Processor grid} = \begin{vmatrix} 0 \\ 1 \\ 2 \\ 3 \end{vmatrix} = \begin{vmatrix} (0,0) \\ (0,1) \\ (1,0) \\ (1,1) \end{vmatrix}$$

DMAT
○○○○○○○○○
○○○○●○○
○○○○○○○○○

pbdDMAT eg's
○○○○
○○○

DMAT Distributions

## DMAT: 2-dimensional Row Cyclic

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

$$\text{Processor grid} = \begin{vmatrix} 0 & 1 \\ 2 & 3 \end{vmatrix} = \begin{vmatrix} (0,0) & (0,1) \\ (1,0) & (1,1) \end{vmatrix}$$

DMAT
○○○○○○○○○
○○○○○○●
○○○○○○○○○

pbdDMAT eg's
○○○○
○○○

DMAT Distributions

## DMAT: 2-dimensional Block-Cyclic

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

$$\text{Processor grid} = \begin{vmatrix} 0 & 1 \\ 2 & 3 \end{vmatrix} = \begin{vmatrix} (0,0) & (0,1) \\ (1,0) & (1,1) \end{vmatrix}$$

DMAT
○○○○○○○○○
○○○○○○
●○○○○○○○○
pbdDMAT

pbdDMAT eg's
○○○○
○○○

The `DMAT` Data Structure

The more complicated the processor grid, the more complicated the distribution.

DMAT
○○○○○○○○
○○○○○○
○○●○○○○○○
pbdDMAT

pbdDMAT eg's
○○○○
○○○

## DMAT: 2-dimensional Block-Cyclic with 6 Processors

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

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

DMAT
○○○○○○○○○
○○○○○○
○○●○○○○○○

pbdDMAT

pbdDMAT eg's
○○○○
○○○

## Understanding DMAT: Local View

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

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

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

DMAT
○○○○○○○○○
○○○○○○
○○○○○●○○○○

pbdDMAT eg's
○○○○
○○○

pbdDMAT

## The `DMAT` 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 processor grid and blocking factor used.

④ `DMAT` is locally column-major and globally, it depends. . .

$$\begin{bmatrix} x_{11} & x_{12} & x_{13} & x_{14} & x_{15} \\ x_{21} & x_{22} & x_{23} & x_{24} & x_{25} \\ x_{31} & x_{32} & x_{33} & x_{34} & x_{35} \\ x_{41} & x_{42} & x_{43} & x_{44} & x_{45} \\ x_{51} & x_{52} & x_{53} & x_{54} & x_{55} \\ x_{61} & x_{62} & x_{63} & x_{64} & x_{65} \\ x_{71} & x_{72} & x_{73} & x_{74} & x_{75} \\ x_{81} & x_{82} & x_{83} & x_{84} & x_{85} \\ x_{91} & x_{92} & x_{93} & x_{94} & x_{95} \end{bmatrix}$$

⑥ `GBD` is a generalization of the one-dimensional block `DMAT` distribution. Otherwise there is no relation.

⑦ `DMAT` is confusing, but very robust.

DMAT
○○○○○○○○○
○○○○○○
○○○○○●○○○

pbdDMAT eg's
○○○○
○○○

pbdDMAT

## Pros and Cons of This Data Structure

### Pros
- Fast for distributed matrix computations

### Cons
- Literally everything else

*This is why we hide most of the distributed details.*

The details are there if you want them (you don't want them).

**DMAT**
○○○○○○○○○
○○○○○○
○○○○○○●○○

**pbdDMAT eg's**
○○○○
○○○

**pbdDMAT**

### Distributed Matrix Methods

**pbdDMAT** has over 100 methods with *identical* syntax to R:

- `` `[` ``, rbind(), cbind(), ...
- lm.fit(), prcomp(), cov(), ...
- `` `%*%` ``, solve(), svd(), norm(), ...
- median(), mean(), rowSums(), ...

Serial Code

```
1  cov(x)
```

Parallel Code

```
1  cov(x)
```

**DMAT**
○○○○○○○○○
○○○○○○
○○○○○○○●○

pbdDMAT

**pbdDMAT eg's**
○○○○
○○○

## Comparing pbdMPI and pbdDMAT

**pbdMPI**:

- MPI + sugar.
- GBD not the only structure **pbdMPI** can handle (just a useful convention).

**pbdDMAT**:

- More of a software package.
- DMAT structure *must* be used for **pbdDMAT**.
- If the data is not 2d block-cyclic compatible, DMAT will *definitely* give the wrong answer.

**DMAT**
○○○○○○○○○
○○○○○○
○○○○○○○●

**pbdDMAT eg's**
○○○○
○○○

**pbdDMAT**

## Quick Comments for Using pbdDMAT

**❶** Start by loading the package:

```
1  library(pbdDMAT, quiet = TRUE)
```

**❷** Always initialize before starting and finalize when finished:

```
1  init.grid()
2
3  # ...
4
5  finalize()
```

**❸** Distributed `DMAT` objects will be given the suffix `.dmat` to visually help distinguish them from global objects. This suffix carries no semantic meaning.

**DMAT**
○○○○○○○○○
○○○○○○
○○○○○○○○

**pbdDMAT eg's**
○○○○
○○○

# Contents

**DMAT**
○○○○○○○○
○○○○○○
○○○○○○○○○

**pbdDMAT eg's**
●○○○
○○○

Statistics Examples with pbdDMAT

## Sample Covariance

### Serial Code

```
1 Cov.X <- cov(X)
2 print(Cov.X)
```

### Parallel Code

```
1 Cov.X <- cov(X)
2 print(Cov.X)
```

DMAT
○○○○○○○○○
○○○○○○
○○○○○○○○○

pbdDMAT eg's
○●○○
○○○

Statistics Examples with pbdDMAT

## Linear Regression

### Serial Code

```
1  tX <- t(X)
2  A <- tX %*% X
3  B <- tX %*% y
4
5  ols <- solve(A) %*% B
6
7  # or
8  ols <- lm.fit(X, y)
```

### Parallel Code

```
1  tX <- t(X)
2  A <- tX %*% X
3  B <- tX %*% y
4
5  ols <- solve(A) %*% B
6
7  # or
8  ols <- lm.fit(X, y)
```

DMAT
○○○○○○○○○
○○○○○○
○○○○○○○○○

pbdDMAT eg's
○○○●○
○○○○
○○○

Statistics Examples with pbdDMAT

## Example 5: PCA

PCA: pca.r

```
1   library(pbdDMAT, quiet=T)
2   init.grid()
3
4   n <- 1e4
5   p <- 250
6
7   comm.set.seed(diff=T)
8   x.dmat <- ddmatrix("rnorm", nrow=n, ncol=p, mean=100, sd=25)
9
10  pca <- prcomp(x=x.dmat, retx=TRUE, scale=TRUE)
11  prop_var <- cumsum(pca$sdev)/sum(pca$sdev)
12  i <- max(min(which(prop_var > 0.9)) - 1, 1)
13
14  y.dmat <- pca$x[, 1:i]
15
16  comm.cat("\nCols: ", i, "\n", quiet=T)
17  comm.cat("%Cols:", i/dim(x.dmat)[2], "\n\n", quiet=T)
18
19  finalize()
```

Execute this script via:

```
1   mpirun -np 2 Rscript 5_pca.r
```

Sample Output:

```
1   Cols:   221
2   %Cols:  0.884
```

**DMAT**
○○○○○○○○○
○○○○○○
○○○○○○○○○

**pbdDMAT eg's**
○○○●
○○○

**Statistics Examples with pbdDMAT**

### Distributed Matrices

**pbdDEMO** contains many other examples of reading and managing GBD and DMAT data

DMAT
○○○○○○○○○
○○○○○○
○○○○○○○○○

pbdDMAT eg's
○○○○
●○○

RandSVD

## Randomized SVD[1]

### PROTOTYPE FOR RANDOMIZED SVD

*Given an $m \times n$ matrix $A$, a target number $k$ of singular vectors, and an exponent $q$ (say, $q = 1$ or $q = 2$), this procedure computes an approximate rank-$2k$ factorization $U\Sigma V^*$, where $U$ and $V$ are orthonormal, and $\Sigma$ is nonnegative and diagonal.*

**Stage A:**
1. Generate an $n \times 2k$ Gaussian test matrix $\Omega$.
2. Form $Y = (AA^*)^q A\Omega$ by multiplying alternately with $A$ and $A^*$.
3. Construct a matrix $Q$ whose columns form an orthonormal basis for the range of $Y$.

**Stage B:**
4. Form $B = Q^*A$.
5. Compute an SVD of the small matrix: $B = \tilde{U}\Sigma V^*$.
6. Set $U = Q\tilde{U}$.

**Note:** The computation of $Y$ in step 2 is vulnerable to round-off errors. When high accuracy is required, we must incorporate an orthonormalization step between each application of $A$ and $A^*$; see Algorithm 4.4.

### ALGORITHM 4.4: RANDOMIZED SUBSPACE ITERATION

*Given an $m \times n$ matrix $A$ and integers $\ell$ and $q$, this algorithm computes an $m \times \ell$ orthonormal matrix $Q$ whose range approximates the range of $A$.*

1. Draw an $n \times \ell$ standard Gaussian matrix $\Omega$.
2. Form $Y_0 = A\Omega$ and compute its QR factorization $Y_0 = Q_0 R_0$.
3. **for** $j = 1, 2, \ldots, q$
4.     Form $\tilde{Y}_j = A^*Q_{j-1}$ and compute its QR factorization $\tilde{Y}_j = \tilde{Q}_j \tilde{R}_j$.
5.     Form $Y_j = A\tilde{Q}_j$ and compute its QR factorization $Y_j = Q_j R_j$.
6. **end**
7. $Q = Q_q$.

### Serial R

```r
1   randSVD <- function(A, k, q=3)
2     {
3       ## Stage A
4       Omega <-  matrix(rnorm(n*2*k),
5                        nrow=n, ncol=2*k)
6       Y <- A %*% Omega
7       Q <- qr.Q(qr(Y))
8       At <- t(A)
9       for(i in 1:q)
10        {
11          Y <- At %*% Q
12          Q <- qr.Q(qr(Y))
13          Y <- A %*% Q
14          Q <- qr.Q(qr(Y))
15        }
16
17      ## Stage B
18      B <- t(Q) %*% A
19      U <- La.svd(B)$u
20      U <- Q %*% U
21      U[, 1:k]
22    }
```

[1] Halko N, Martinsson P-G and Tropp J A 2011 Finding structure with randomness: probabilistic algorithms for constructing approximate matrix decompositions *SIAM Rev.* **53** 217–88

DMAT
○○○○○○○○○
○○○○○○
○○○○○○○○○

pbdDMAT eg's
○○○○
○●○
RandSVD

## Randomized SVD

### Serial R

```
 1  randSVD <- function(A, k, q=3)
 2    {
 3      ## Stage A
 4      Omega <-  matrix(rnorm(n*2*k),
 5              nrow=n, ncol=2*k)
 6      Y <- A %*% Omega
 7      Q <- qr.Q(qr(Y))
 8      At <- t(A)
 9      for(i in 1:q)
10        {
11          Y <- At %*% Q
12          Q <- qr.Q(qr(Y))
13          Y <- A %*% Q
14          Q <- qr.Q(qr(Y))
15        }
16
17      ## Stage B
18      B <- t(Q) %*% A
19      U <- La.svd(B)$u
20      U <- Q %*% U
21      U[, 1:k]
22    }
```

### Parallel pbdR
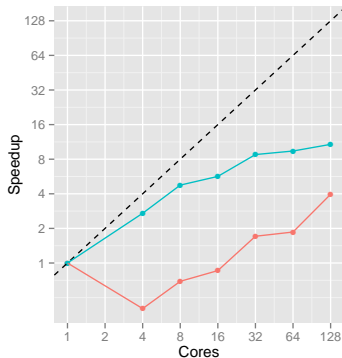
```
 1  randSVD <- function(A, k, q=3)
 2    {
 3      ## Stage A
 4      Omega <-  ddmatrix("rnorm",
 5              nrow=n, ncol=2*k)
 6      Y <- A %*% Omega
 7      Q <- qr.Q(qr(Y))
 8      At <- t(A)
 9      for(i in 1:q)
10        {
11          Y <- At %*% Q
12          Q <- qr.Q(qr(Y))
13          Y <- A %*% Q
14          Q <- qr.Q(qr(Y))
15        }
16
17      ## Stage B
18      B <- t(Q) %*% A
19      U <- La.svd(B)$u
20      U <- Q %*% U
21      U[, 1:k]
22    }
```

DMAT
○○○○○○○○○
○○○○○○
○○○○○○○○○

pbdDMAT eg's
○○○○
○○●

RandSVD

## Randomized SVD



30 Singular Vectors from a 100,000 by 1,000 Matrix

30 Singular Vectors from a 100,000 by 1,000 Matrix
Speedup of Randomized vs. Full SVD