

# From 1 Core to Thousands: R to pbdR

Drew Schmidt<sup>1</sup>, George Ostrouchov<sup>2</sup>, Wei-Chen Chen<sup>2</sup>,  
Pragneshkumar Patel<sup>1</sup>

1. University of Tennessee, USA
2. Oak Ridge National Laboratory, USA

August 8, 2013



# Affiliations and Support

The pbdR Core Team

<http://r-pbd.org>

Wei-Chen Chen<sup>1</sup>, George Ostrouchov<sup>1,2</sup>, Pragneshkumar Patel<sup>2</sup>, Drew Schmidt<sup>1</sup>

Ostrouchov, Patel, and Schmidt 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.

---

<sup>1</sup>Computer Science and Mathematics Division, Oak Ridge National Laboratory, Oak Ridge, TN

<sup>2</sup>Remote Data Analysis and Visualization Center, University of Tennessee, Knoxville, TN

# 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 Nautilus from:  
[/lustre/medusa/mschmid3/tutorial/scripts.tar.gz](http://lustre/medusa/mschmid3/tutorial/scripts.tar.gz)

# Contents

- 1 Introduction to R
- 2 Basic R Syntax
- 3 pbdR
- 4 Benchmarks
- 5 Challenges

# Contents

- 1 Introduction to R
  - What is R?
  - Syntax for Data Science

## What is R?

- *lingua franca* for data analytics and statistical computing.
- Part programming language, part data analysis package.
- Dialect of S (Bell Labs).
- Syntax designed for data. scoping semantics, and 2 official OOP systems.

## Who uses R?

Google, Pfizer, Merck, Bank of America, Shell<sup>a</sup>,  
Oracle<sup>b</sup>, Facebook, bing, Mozilla, okcupid<sup>c</sup>,  
ebay<sup>d</sup>,  
kickstarter<sup>e</sup>, the New York Times<sup>f</sup>

<sup>a</sup>[https://www.nytimes.com/2009/01/07/technology/business-computing/07program.html?\\_r=0](https://www.nytimes.com/2009/01/07/technology/business-computing/07program.html?_r=0)

<sup>b</sup><http://www.oracle.com/us/corporate/features/features-oracle-r-enterprise-498732.html>

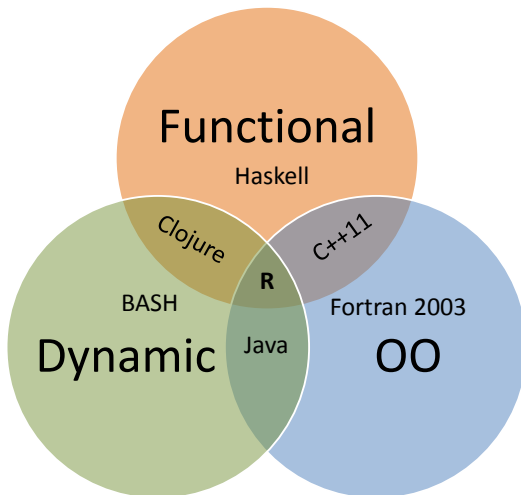
<sup>c</sup><http://www.revolutionanalytics.com/what-is-open-source-r/companies-using-r.php>

<sup>d</sup><http://blog.revolutionanalytics.com/2012/09/using-r-in-production-industry-experts-share-their-experiences.html>

<sup>e</sup><http://blog.revolutionanalytics.com/2012/09/kickstarter-facilitates-50m-in-indie-game-funding.html>

<sup>f</sup><http://blog.revolutionanalytics.com/2012/05/nyt-charts-the-facebook-ipo-with-r.html>

## Language Paradigms





## Data Types

- Storage: logical, int, double, double complex, character
- Structures: vector, matrix, array, list, dataframe
- Caveats: (Logical) TRUE, FALSE, NA

For the remainder of the tutorial, we will restrict ourselves to real number matrix computations.

## High Level Syntax

```

1 x <- matrix(rnorm(30), nrow=10)
2 x <- x[-1, 2:5]
3 x <- log(abs(x) + 1)
4 xtx <- t(x) %*% x
5 ans <- svd(solve(xtx))

```

## More than just a Matlab clone. . .

- Data science (machine learning, statistics, data mining, . . . ) is mostly matrix algebra.

So what about Matlab/Python/Julia/. . . ?

- Depends on your “religion”
- As a *data analysis* package, R is king.

## High Level Syntax *for Data*

```
1 pca <- prcomp(x, retx=TRUE, scale=TRUE)
2 prop_var <- cumsum(pca$sdev)/sum(pca$sdev)
3 i <- min(which(prop_var > 0.9)) - 1
4
5 y <- pca$x[, 1:i]
```

## Basics (1 of 2)

- The default method is to print:

```
1 R> sum
2 function (... , na.rm = FALSE) .Primitive("sum")
```

- Use <- for assignment:

```
1 R> x <- 1
2 R> x+1
3 [1] 2
```

- Naming rules: mostly like C.
- R is case sensitive.
- We use . the way most languages use \_, e.g., La.svd() instead of La\_svd().
- We use \$ (sometimes @) the way most languages use .

## Basics (2 of 2)

- Use ? or ?? to search help

```
1 R> ?set.seed
2 R> ?comm.set.seed
3 No documentation for comm.set.seed in
  specified packages and libraries:
4 you could try ??comm.set.seed
5 R> ??comm.set.seed
```

## Addons and Extras

R has the Comprehensive R Archive Network (CRAN), which is a package repository like CTAN and CPAN.

From R

```
1 install.packages("pbdMPI") # install
2 library(pbdMPI)           # load
```

From Shell

```
1 R CMD INSTALL pbdMPI_0.1-6.tar.gz
```

## Lists (1 of 1)

```

1 R> l <- list(a=1, b="a")
2 R> l
3 $a
4 [1] 1
5
6 $b
7 [1] "a"
8
9 R> l$a
10 [1] 1
11
12 R> list(x=list(a=1, b="a"), y=TRUE)
13 $x
14 $x$a
15 [1] 1
16
17 $x$b
18 [1] "a"
19
20
21 $y
22 [1] TRUE

```



## Vectors and Matrices (1 of 2)

```

1 R> c(1, 2, 3, 4, 5, 6)
2 [1] 1 2 3 4 5 6
3
4 R> matrix(1:6, nrow=2, ncol=3)
5      [,1] [,2] [,3]
6 [1,]    1    3    5
7 [2,]    2    4    6
8
9 R> x <- matrix(1:6, nrow=2, ncol=3)
10
11 R> x[, -1]
12      [,1] [,2]
13 [1,]    3    5
14 [2,]    4    6
15
16 R> x[1, 1:2]
17 [1] 1 3

```

## Vectors and Matrices (2 of 2)

```

1 R> dim(x)
2 [1] 2 3
3
4 R> dim(x) <- NULL
5 R> x
6 [1] 1 2 3 4 5 6
7
8 R> dim(x) <- c(3,2)
9 R> x
10      [,1] [,2]
11 [1,]    1    4
12 [2,]    2    5
13 [3,]    3    6

```

## Vector and Matrix Arithmetic (1 of 2)

```

1 R> 1:4 + 4:1
2 [1] 5 5 5 5
3
4 R> x <- matrix(0, nrow=2, ncol=3)
5
6 R> x + 1
7      [,1] [,2] [,3]
8 [1,]    1    1    1
9 [2,]    1    1    1
10
11 R> x + 1:3
12      [,1] [,2] [,3]
13 [1,]    1    3    2
14 [2,]    2    1    3

```

## Vector and Matrix Arithmetic (2 of 2)

```

1 R> x <- matrix(1:6, nrow=2)
2
3 R> x*x
4      [,1] [,2] [,3]
5 [1,]    1    9   25
6 [2,]    4   16   36
7
8 R> x %*% x
9 Error in x %*% x : non-conformable arguments
10
11 R> t(x) %*% x
12      [,1] [,2] [,3]
13 [1,]    5   11   17
14 [2,]   11   25   39
15 [3,]   17   39   61
16
17 R> crossprod(x)
18      [,1] [,2] [,3]
19 [1,]    5   11   17
20 [2,]   11   25   39
21 [3,]   17   39   61

```

## Linear Algebra (1 of 2): Matrix Inverse

$$x_{n \times n} \text{ invertible} \iff \exists y_{n \times n} (xy = yx = Id_{n \times n})$$

```

1 R> x <- matrix(rnorm(5*5), nrow=5)
2 R> y <- solve(x)
3
4 R> round(x %*% y)
5      [,1] [,2] [,3] [,4] [,5]
6 [1,]    1    0    0    0    0
7 [2,]    0    1    0    0    0
8 [3,]    0    0    1    0    0
9 [4,]    0    0    0    1    0
10 [5,]    0    0    0    0    1

```

## Linear Algebra (2 of 2): Singular Value Decomposition

$$x = U\Sigma V^T$$

```

1 R> x <- matrix(rnorm(2*3), nrow=3)
2 R> svd(x)
3 $d
4 [1] 2.4050716 0.3105008
5
6 $u
7           [,1]      [,2]
8 [1,] 0.8582569 -0.1701879
9 [2,] 0.2885390  0.9402076
10 [3,] 0.4244295 -0.2950353
11
12 $v
13           [,1]      [,2]
14 [1,] -0.05024326 -0.99873701
15 [2,] -0.99873701  0.05024326

```

## More than just a Matlab clone. . .

- Data science (machine learning, statistics, data mining, . . . ) is mostly matrix algebra.

So what about Matlab/Python/Julia/. . . ?

- The one you prefer depends more on your “religion” rather than differences in capabilities.
- As a *data analysis* package, R is king.

## Simple Statistics (1 of 2): Summary Statistics

```
1 R> x <- matrix(rnorm(30, mean=10, sd=3), nrow=10)
2
3 R> mean(x)
4 [1] 9.825177
5
6 R> median(x)
7 [1] 9.919243
8
9 R> sd(as.vector(x))
10 [1] 3.239388
11
12 R> colMeans(x)
13 [1] 9.661822 10.654686 9.159025
14
15 R> apply(x, MARGIN=2, FUN=sd)
16 [1] 2.101059 3.377347 4.087131
```



## Simple Statistics (2 of 2): Sample Covariance

$$\text{cov}(x_{n \times p}) = \frac{1}{n-1} \sum_{i=1}^n (x_i - \mu_x)(x_i - \mu_x)^T$$

```
1 x <- matrix(rnorm(30), nrow=10)
2
3 # least recommended
4 cm <- colMeans(x)
5 crossprod(sweep(x, MARGIN=2, STATS=cm))
6
7 # less recommended
8 crossprod(scale(x, center=TRUE, scale=FALSE))
9
10 # recommended
11 cov(x)
```

## Advanced Statistics (1 of 2): Principal Components

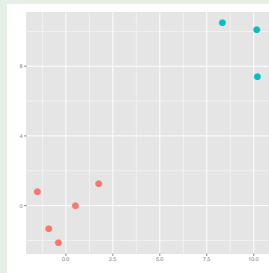
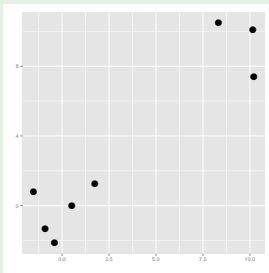
PCA = centering + scaling + rotation (via SVD)

```

1 R> x <- matrix(rnorm(30), nrow=10)
2
3 R> prcomp(x, retx=TRUE, scale=TRUE)
4 Standard deviations:
5 [1] 1.1203373 1.0617440 0.7858397
6
7 Rotation:
8
9           PC1          PC2          PC3
10 [1,]  0.71697825 -0.3275365  0.6153552
11 [2,] -0.03382385  0.8653562  0.5000147
12 [3,]  0.69627447  0.3793133 -0.6093630

```

## Advanced Statistics (2 of 2): k-Means Clustering



```
1 R> x <- rbind(matrix(rnorm(5*2, mean=0), ncol=2),
2               matrix(rnorm(3*2, mean=10), ncol=2))
```

## Advanced Statistics (2 of 2): k-Means Clustering

```

1 R> kmeans(x, centers=2)
2 K-means clustering with 2 clusters of sizes 5, 3
3
4 Cluster means:
5      [,1]      [,2]
6 1 -0.1080612 -0.2827576
7 2  9.5695365  9.3191892
8
9 Clustering vector:
10 [1] 1 1 1 1 1 2 2 2
11
12 Within cluster sum of squares by cluster:
13 [1] 14.675072  7.912641
14 (between_SS / total_SS =  93.9 %)
15
16 Available components:
17
18 [1] "cluster"      "centers"      "totss"
19 [2] "withinss"     "tot.withinss"
20 [3] "betweenss"    "size"

```

# Contents

## 3 pbdR

- The pbdR Project
- pbdR Paradigms

## Programming with Big Data in R (pbdR)

Striving for *Productivity, Portability, Performance*

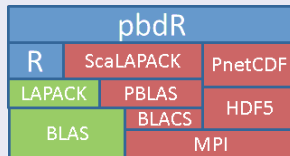
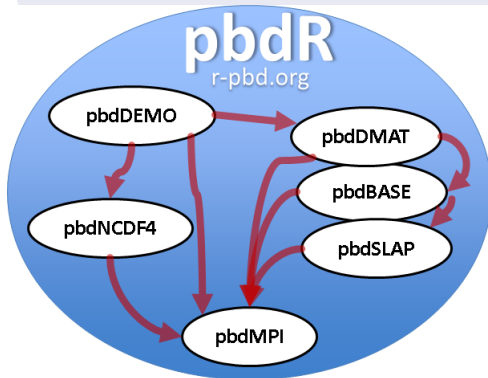


- *Free<sup>a</sup>* R packages.
- Bridging high-performance C with high-productivity of R
- Scalable, big data analytics.
- Distributed data details implicitly managed.
- Methods have syntax *identical* to R.
- Powered by state of the art numerical libraries (MPI, ScaLAPACK, ...)

---

<sup>a</sup>MPL, BSD, and GPL licensed

## pbdR Packages



## pbdR on HPC Resources

pbdR is currently installed and maintained on:

- Nautilus, UTK
- Kraken, UTK
- Newton, UTK
- Lens, ORNL
- Titan, ORNL
- tara, UMBC

If you are interested in maintaining pbdR, contact us at

[RBigData@gmail.com](mailto:RBigData@gmail.com)



## Example Syntax

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

Look familiar?

*The above runs on 1 core with R or 10,000 cores with pbdR*

## pbdR Paradigms

Programs that use pbdR utilize:

- Batch execution
- Single Program/Multiple Data (SPMD) style

And generally utilize:

- Data Parallelism

## Batch Execution

- Non-interactive
- Use

```
1 Rscript my_script.r
```

or

```
1 R CMD BATCH my_script.r
```

- In parallel:

```
1 mpirun -np 2 Rscript my_par_script.r
```

## Single Program/Multiple Data (SPMD)

- Difficult to describe, easy to do. . .
- Only one program is written, executed in batch on all processors.
- Different processors are autonomous; there is no manager.
- The dominant programming model for large machines.

# Contents

1 Introduction to R

2 Basic R Syntax

3 pbdR

4 Benchmarks

5 Challenges

## Non-Optimal Choices Throughout

- 1 Only libre software used (no MKL, ACML, etc.).
- 2 1 core = 1 MPI process.
- 3 No tuning for data distribution.

## Benchmark Data

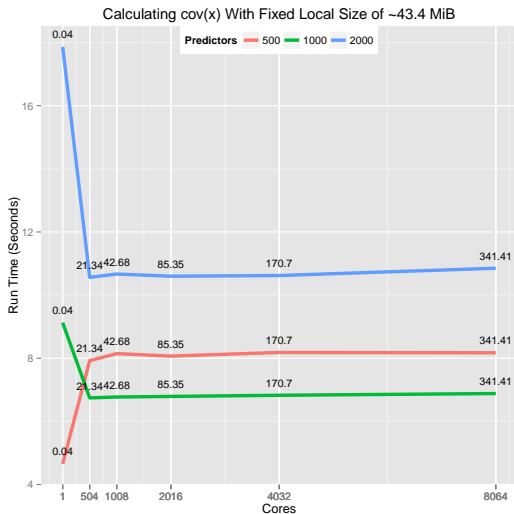
- 1 Random normal  $N(100, 10000)$ .
- 2 Local problem size of  $\approx 43.4 \text{ MiB}$ .
- 3 Three sets: 500, 1000, and 2000 columns.
- 4 Several runs at different core sizes within each set.

## Covariance Code

```
1 x <- ddmatrix("rnorm", nrow=n, ncol=p, mean=mean, sd=sd)
2
3 cov.x <- cov(x)
```



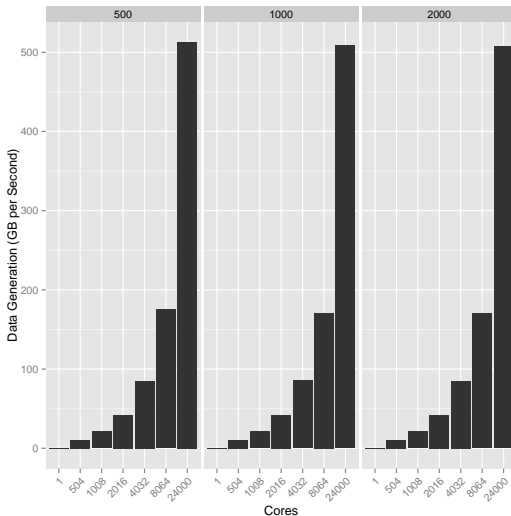
cov()



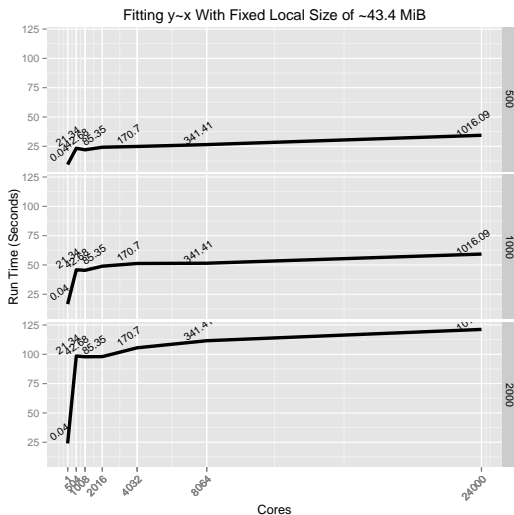
## Linear Model Code

```
1 x <- ddmatrix("rnorm", nrow=n, ncol=p, mean=mean, sd=sd)
2 beta_true <- ddmatrix("runif", nrow=p, ncol=1)
3
4 y <- x %*% beta_true
5
6 beta_est <- lm.fit(x=x, y=y)$coefficients
```

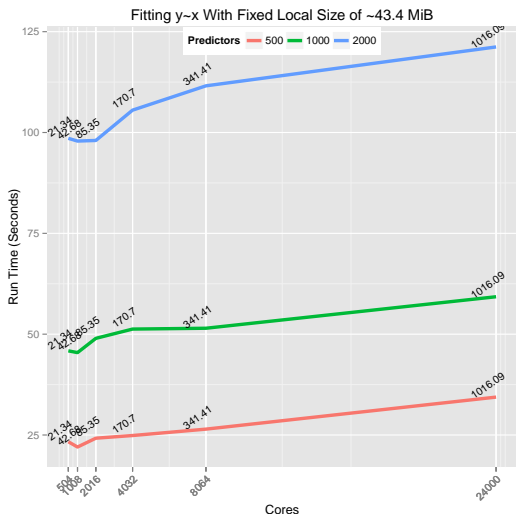
## Data Generation



lm.fit()



```
lm.fit()
```



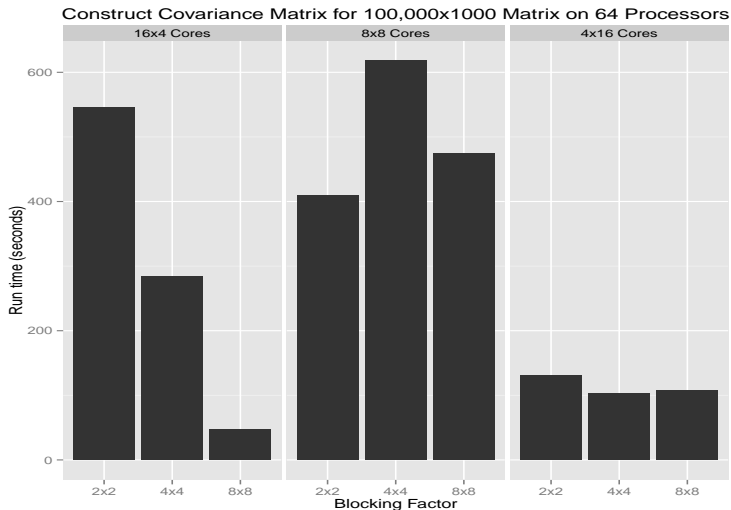
# Contents

- 1 Introduction to R
- 2 Basic R Syntax
- 3 pbdR
- 4 Benchmarks
- 5 Challenges**

## Challenges

- Perceptions.
- Library loading.
- Profiling.

## Covariance Revisited: Distributed Data Parameter Calibration





## Tutorials

- SC13, November 17-22, Denver, Colorado, USA

## Invited Talks

- IASC, Aug 22-23, Seoul
- World Statistics Congress, August 25-30, Hong Kong

Thanks for coming!

Questions?

Be sure to stick around for the tutorial