

# John Myles White

"Who refuses to do arithmetic is doomed to talk nonsense."

*Browse: [Home](#) / [2011](#) / [October](#) / [31](#) / [Using Sparse Matrices in R](#)*

## Using Sparse Matrices in R

*By [John Myles White](#) on 10.31.2011*

### Introduction

I've recently been working with a couple of large, extremely sparse data sets in R. This has pushed me to spend some time trying to master the CRAN packages that support sparse matrices. This post describes three of them: the `Matrix`, `slam` and `glmnet` packages. The first two packages provide data storage classes for sparse matrices, while the last package can perform GLM analyses on data stored in a sparse matrix.

### The Matrix Package

The first package I worked with that provides a sparse matrix implementation is Doug Bates' [Matrix](#) package. In fact, matrices of class `Matrix` can be switched between full and sparse representations dynamically, but I'll focus on forcing the use of a sparse representation.

As a first example, it's helpful to generate a 1000×1000 matrix of zeros using the `matrix` class and then another 1000×1000 matrix of zeros using the `Matrix` class:

```
1 library('Matrix')
2
3 m1 <- matrix(0, nrow = 1000, ncol = 1000)
4 m2 <- Matrix(0, nrow = 1000, ncol = 1000, sparse = TRUE)
5
6 object.size(m1)
7 # 8000200 bytes
8 object.size(m2)
9 # 5632 bytes
```

As you can see, the representation of a perfectly sparse matrix using the `Matrix` class uses more than one thousand times as much RAM. And you can check how much more RAM would be required if both matrices had exactly one non-zero entry:

```
1 m1[500, 500] <- 1
2 m2[500, 500] <- 1
3
4 object.size(m1)
5 # 8000200 bytes
6 object.size(m2)
7 # 5648 bytes
```

The full matrix representation does not change in size because all of the zeros are being represented explicitly, while the sparse matrix is conserving that space by representing only the non-zero entries. With a simple subtraction, we find that adding one additional non-zero entry increases the storage requirement by 16 bytes. We can conclude that setting all of the entries to non-zero values would require  $5632 + 16 \times$

1000 \* 1000 bytes, which is 16,005,632 bytes or almost exactly twice the amount of storage required to use the full representation implemented by the `matrix` class.

Of course, the take away lesson is that sparse matrices are very efficient if your data is sparse and mildly wasteful if your data is not sparse.

Beyond simple initialization and assignment operations, we can perform quite a few other operations on objects of class `Matrix`, including vector multiplication, matrix addition and subtraction and transposition:

1 2 3 4	<pre>m2 %*% rnorm(1000) m2 + m2 m2 - m2 t(m2)</pre>
------------------	---

In addition you can do binding operations using objects of class `Matrix` with the `cBind` and `rBind` functions:

1 2 3 4 5 6 7	<pre>m3 &lt;- cBind(m2, m2) nrow(m3) ncol(m3)  m4 &lt;- rBind(m2, m2) nrow(m4) ncol(m4)</pre>
---------------------------------	---

Other operations are probably supported, but I haven't need them so far in my work.

## The slam Package

An alternative to the `Matrix` package is the `slam` package by Kurt Hornik and others. The sparse matrices generated using this package can be noticeably smaller than those generated by the `Matrix` package in some cases.

For example, the same perfectly sparse matrix using the `slam` package requires only 1,032 bytes of space:

1 2 3 4 5 6 7 8 9	<pre>library('slam')  m1 &lt;- matrix(0, nrow = 1000, ncol = 1000) m2 &lt;- simple_triplet_zero_matrix(nrow = 1000, ncol = 1000)  object.size(m1) # 8000200 bytes object.size(m2) # 1032 bytes</pre>
---	--

In short, you can make your data fit in 1/5 the amount of RAM under some circumstances. And we can once again check how much additional RAM you need for every new non-zero entry:

1 2 3 4 5 6	<pre># BUG HERE m1[500, 500] &lt;- 1 m2[500, 500] &lt;- 1  object.size(m1) object.size(m2)</pre>
----------------------------	--

Finally, all of the same matrix operations available with the `Matrix` class work on objects implemented by `slam`:

```
1
2
3
4
m2 %*% rnorm(1000)
m2 + m2
m2 - m2
t(m2)
```

## The glmnet Package

Of course, being able to load sparse data into RAM is only interesting if we can analyze it statistically. For me, this usually means that I fit some sort of GLM to the data: most of the time either linear or logistic regression — preferably with some sort of regularization.

Thankfully, the `glmnet` package allows full and sparse matrices to be used without any code changes. You simply need to provide arguments that use the `Matrix` package's `Matrix` class to have `glmnet` switch over to computing with sparse matrices:

```
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
library('Matrix')
library('glmnet')

n <- 10000
p <- 500

x <- matrix(rnorm(n * p), n, p)
iz <- sample(1:(n * p),
             size = n * p * 0.85,
             replace = FALSE)
x[iz] <- 0

object.size(x)

sx <- Matrix(x, sparse = TRUE)

object.size(sx)

beta <- rnorm(p)

y <- x %*% beta + rnorm(n)

glmnet.fit <- glmnet(x, y)
```

How much more efficient is the use of sparse matrices? Here I've recreated the example given in the `glmnet` package docs and added some additional data collection to give a sense of the time and speed gains that come with using sparse matrices for sparse data:

```
1
2
3
4
5
6
7
8
9
10
11
12
library('Matrix')
library('glmnet')

set.seed(1)
performance <- data.frame()

for (sim in 1:10)
{
  n <- 10000
  p <- 500

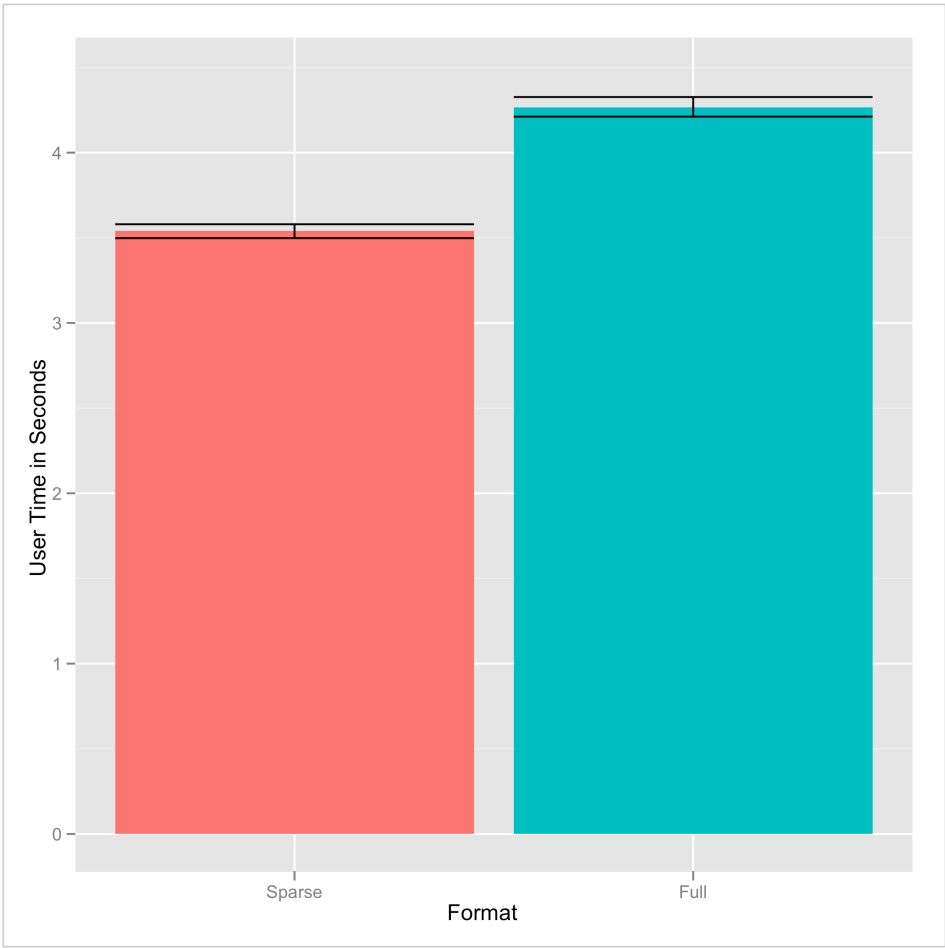
  nzc <- trunc(p / 10)
```

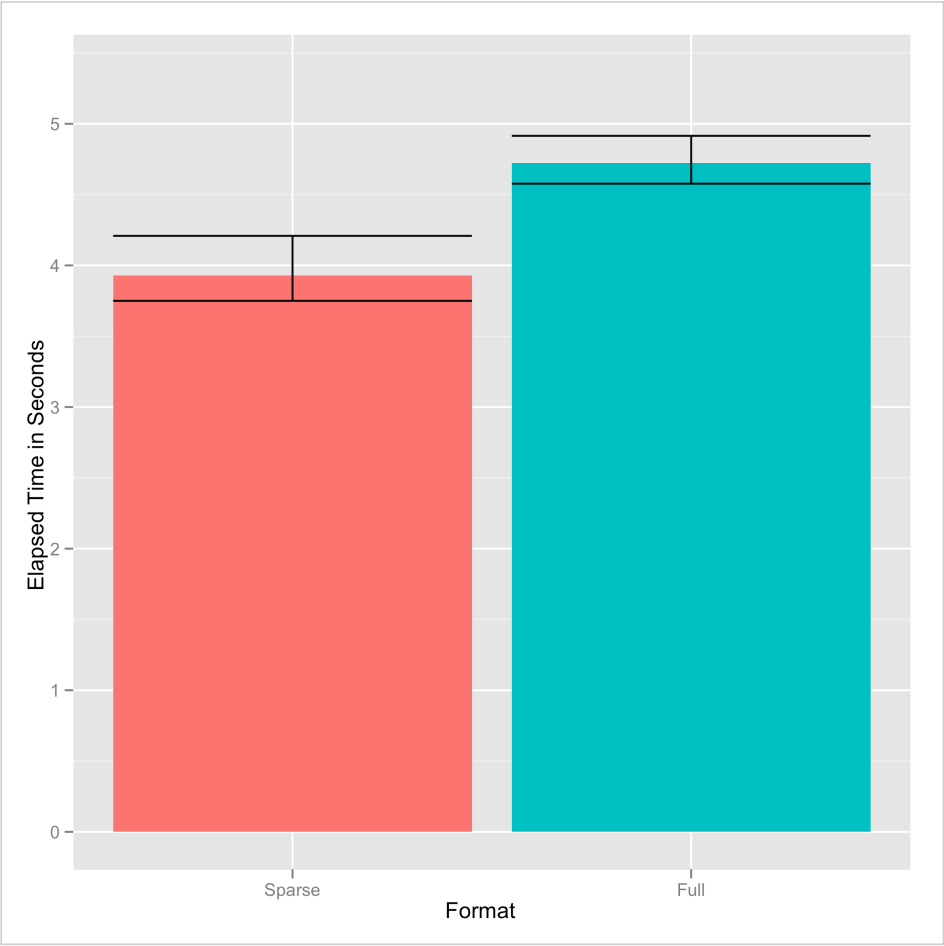
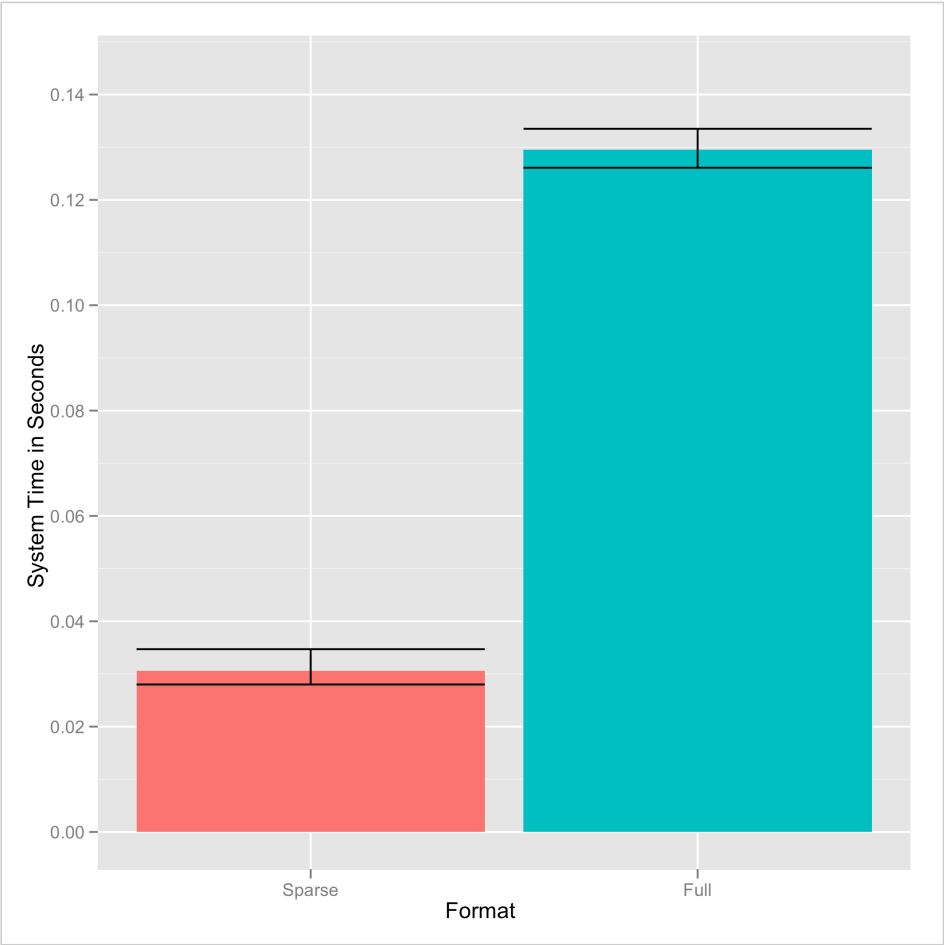
```

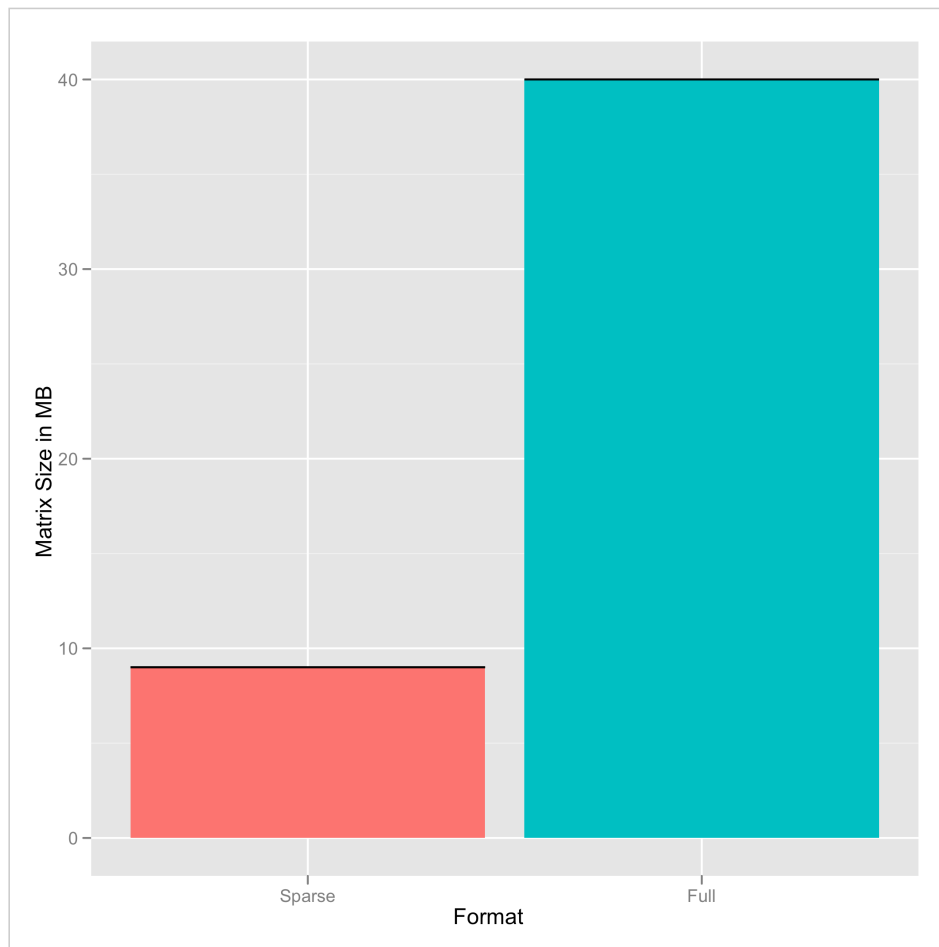
13
14 x <- matrix(rnorm(n * p), n, p)
15 iz <- sample(1:(n * p),
16             size = n * p * 0.85,
17             replace = FALSE)
18 x[iz] <- 0
19 sx <- Matrix(x, sparse = TRUE)
20
21 beta <- rnorm(nzc)
22 fx <- x[, seq(nzc)] %*% beta
23
24 eps <- rnorm(n)
25 y <- fx + eps
26
27 sparse.times <- system.time(fit1 <- glmnet(sx, y))
28 full.times <- system.time(fit2 <- glmnet(x, y))
29 sparse.size <- as.numeric(object.size(sx))
30 full.size <- as.numeric(object.size(x))
31
32 performance <- rbind(performance, data.frame(Format = 'Sparse',
33                                             UserTime = sparse.times[1],
34                                             SystemTime = sparse.times[2],
35                                             ElapsedTime = sparse.times[3],
36                                             Size = sparse.size))
37 performance <- rbind(performance, data.frame(Format = 'Full',
38                                             UserTime = full.times[1],
39                                             SystemTime = full.times[2],
40                                             ElapsedTime = full.times[3],
41                                             Size = full.size))
42 }
43
44 ggplot(performance, aes(x = Format, y = UserTime, fill = Format)) +
45   stat_summary(fun.data = 'mean_cl_boot', geom = 'bar') +
46   stat_summary(fun.data = 'mean_cl_boot', geom = 'errorbar') +
47   ylab('User Time in Seconds') +
48   opts(legend.position = 'none')
49 ggsave('sparse_vs_full_user_time.pdf')
50
51 ggplot(performance, aes(x = Format, y = SystemTime, fill = Format)) +
52   stat_summary(fun.data = 'mean_cl_boot', geom = 'bar') +
53   stat_summary(fun.data = 'mean_cl_boot', geom = 'errorbar') +
54   ylab('System Time in Seconds') +
55   opts(legend.position = 'none')
56 ggsave('sparse_vs_full_system_time.pdf')
57
58 ggplot(performance, aes(x = Format, y = ElapsedTime, fill = Format)) +
59   stat_summary(fun.data = 'mean_cl_boot', geom = 'bar') +
60   stat_summary(fun.data = 'mean_cl_boot', geom = 'errorbar') +
61   ylab('Elapsed Time in Seconds') +
62   opts(legend.position = 'none')
63 ggsave('sparse_vs_full_elapsed_time.pdf')
64
65 ggplot(performance, aes(x = Format, y = Size / 1000000, fill = Format)) +
66   stat_summary(fun.data = 'mean_cl_boot', geom = 'bar') +
67   stat_summary(fun.data = 'mean_cl_boot', geom = 'errorbar') +
68   ylab('Matrix Size in MB') +
69   opts(legend.position = 'none')
70 ggsave('sparse_vs_full_memory.pdf')

```

The images below show the resulting performance metrics:







## Questions

I have some remaining questions about the use of sparse matrices that others might be able to answer:

- How can I perform assignment in place on a matrix generated by the slam package?
- How can I convert between sparse matrices generated by the Matrix and slam packages?

Posted in [Programming, Statistics](#) | [2 Responses](#)

## 2 responses to “Using Sparse Matrices in R”



**Brian Dalesandro** 10.5.2012 at 1:24 pm | [Permalink](#)

Thanks for the examples and analysis. Do you have any experience reading sparse matrices into R? The Matrix package has a few read functions that seem to depend on some very specific sparse matrix representations that I've never heard of (readHB & readMM). Do you know of a more general format and function that can be used? The rest seems straightforward.



**John Myles White** 10.6.2012 at 3:35 pm | [Permalink](#)

I'm not aware of any better reading tools in R, although I imagine they exist. In the past I gave up on sparse matrix support in R and switched to Matlab for one project.

