

STAT-6494 Advanced Statistical Computing with R

Homework 3

Wenjie Wang

12 March 2016

1 Exercise: Linear Model with Sparse Model Matrix

Write **R** functions for linear regression with two inputs, model matrix **x** and response vector **y**. One function **lmQr** uses the QR decomposition of x and other **lmChol** uses the Cholesky decomposition of $x'x$. Each function returns a list with three components: **coefficients** for estimated regression coefficients; **stderr** for standard error of the estimates; and **df.residuals** for the degrees of freedom of the residuals. Compare the performance with **lm.fit** from **R** for both dense matrix and sparse matrix of **x** with various dimensions and sparsity.

1.1 Function **lmQr** and **lmChol**

The function **lmQr** is designed using QR decomposition on design matrix **X** to fit linear regression model. It takes response vector **Y** and design matrix **X**, which can be sparse matrices object in **dgCMatrix** class produced by package **Matrix**. And it returns **coefficients** for estimated coefficients, **stderr** for standard error of the estimates, and **df.residuals** for the degrees of freedom of the residuals. The function **lmChol** is similarly designed but using Cholesky decomposition.

The straightforward implementation is done completely in R. Note that for **lmQr** and **lmChol**, the design matrix **X** input can be either dense matrix class object or sparse matrix class object, **dgCMatrix** class object. The corresponding method will be called depending on the object class.

```
library(Matrix)
lmQr <- function (X, y) {
  resQR <- qr(X)
  betaEst <- qr.coef(resQR, y)
  df <- nrow(X) - (Xp <- ncol(X))
  if (is.matrix(X)) {
    matR <- qr.R(resQR)
    xtxInv <- backsolve(matR, backsolve(matR, diag(Xp),
                                         transpose = TRUE))
  } else {
    matR <- qrR(resQR)
    xtxInv <- solve(crossprod(X))
  }
  sigma2Est <- crossprod(y - X %*% betaEst) / df
  stderr <- sqrt(as.numeric(sigma2Est) * diag(xtxInv))
  list(coefficients = betaEst, stderr = stderr,
        df.residuals = df)
}
lmChol <- function (X, y) {
  resChol <- chol(crossprod(X))
  RHS <- crossprod(X, y)
  xtxInv <- chol2inv(resChol)
```

```

betaEst <- xtxInv %*% RHS
df <- nrow(X) - ncol(X)
sigma2Est <- crossprod(y - X %*% betaEst) / df
stderr <- sqrt(as.numeric(sigma2Est) * diag(xtxInv))
list(coefficients = betaEst, stderr = stderr,
      df.residuals = df)
}

```

1.2 Performance Comparison on Simulated Datasets

We generate simulated datasets from linear regression model with various dimension and sparsity. Then we compare the performance of `lm.fit`, `lmQr`, and `lmChol` on those datasets, respectively.

1.2.1 Dense matrix of small dimension

The response Y is set as a vector with length 500. The design matrix is a 500 by 10 matrix. We compare the computing performance of `lmQr`, `lmChol` with `lm.fit` from package `stats` with the help of package `microbenchmark` over the simulated Y and X as follows:

```

set.seed(1216)
n <- 500
p <- 10
X <- matrix(rnorm(n * p), n)
y <- rowSums(X) + rnorm(n)
library(microbenchmark)
ggplot2::autoplot(microbenchmark(lm.fit(X, y), lmQr(X, y),
                                lmChol(X, y), times = 1000))

```

From the comparison shown in Figure 1, we find that `lmChol` is the fastest one. `lm.fit` is a little bit slower than `lmChol` but faster than `lmQr`.

1.2.2 Sparse matrix of small dimension

The sparsity is set to be 0.9.

```

r <- 0.9
X[sample(n * p, floor(n * p * r))] <- 0
Xsp <- methods::as(X, "dgCMatrix")
ggplot2::autoplot(microbenchmark(lm.fit(X, y), lmQr(X, y), lmQr(Xsp, y),
                                lmChol(X, y), lmChol(Xsp, y), times = 100))

```

From the comparison shown in Figure 2, we find that `lmChol` on dense matrix objects is still the fastest one. `lm.fit` is a little bit slower than `lmChol` but faster than `lmQr` on dense matrix objects. Both `lmChol` and `lmQr` run much slower when the object input is a sparse matrix class object instead of a dense matrix object. It indicates that sparse matrix objects and their corresponding methods cannot bring much help when the dimension of the sparse matrix is relatively small.

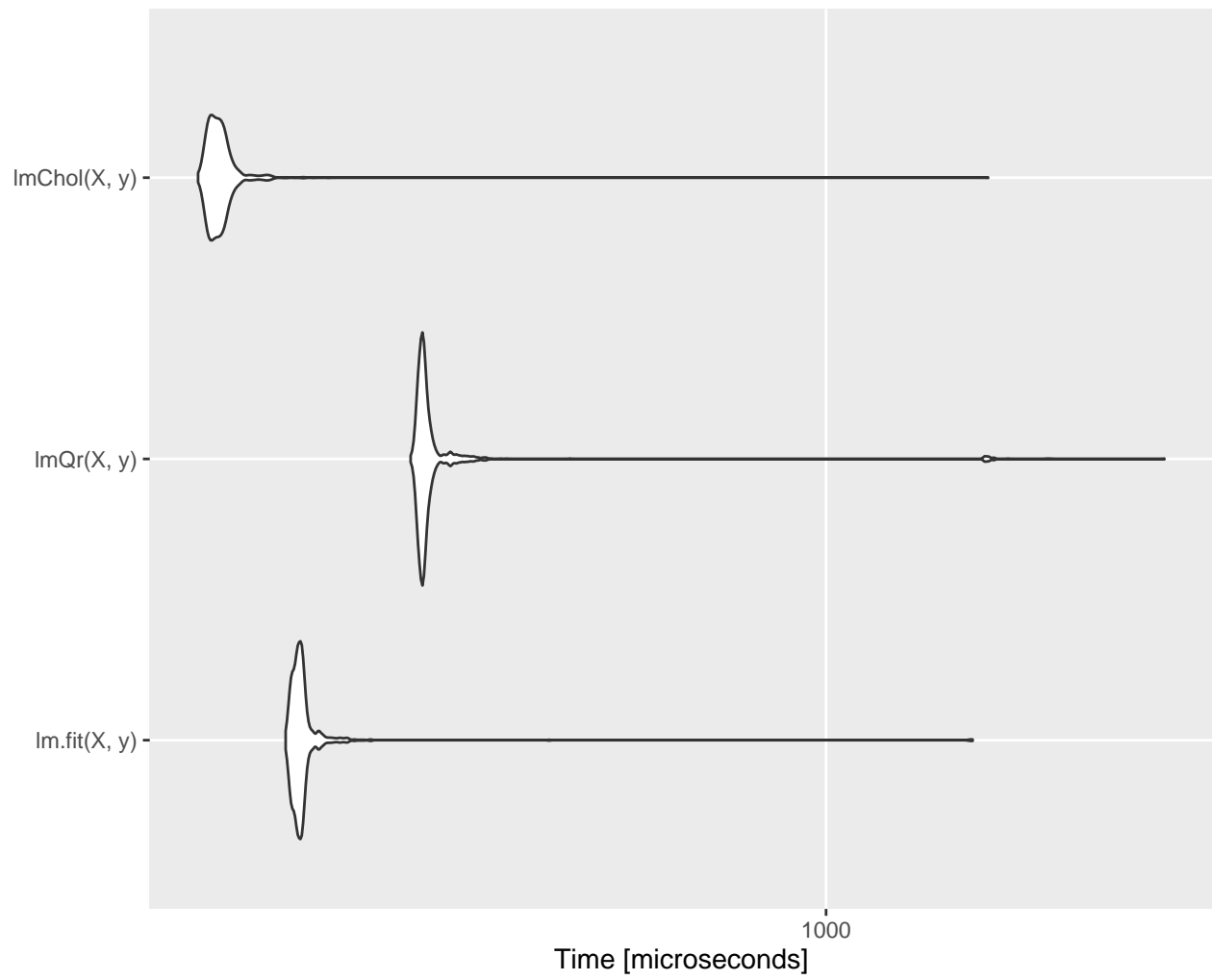


Figure 1: Performance comparison when design matrix is a dense matrix of small dimension

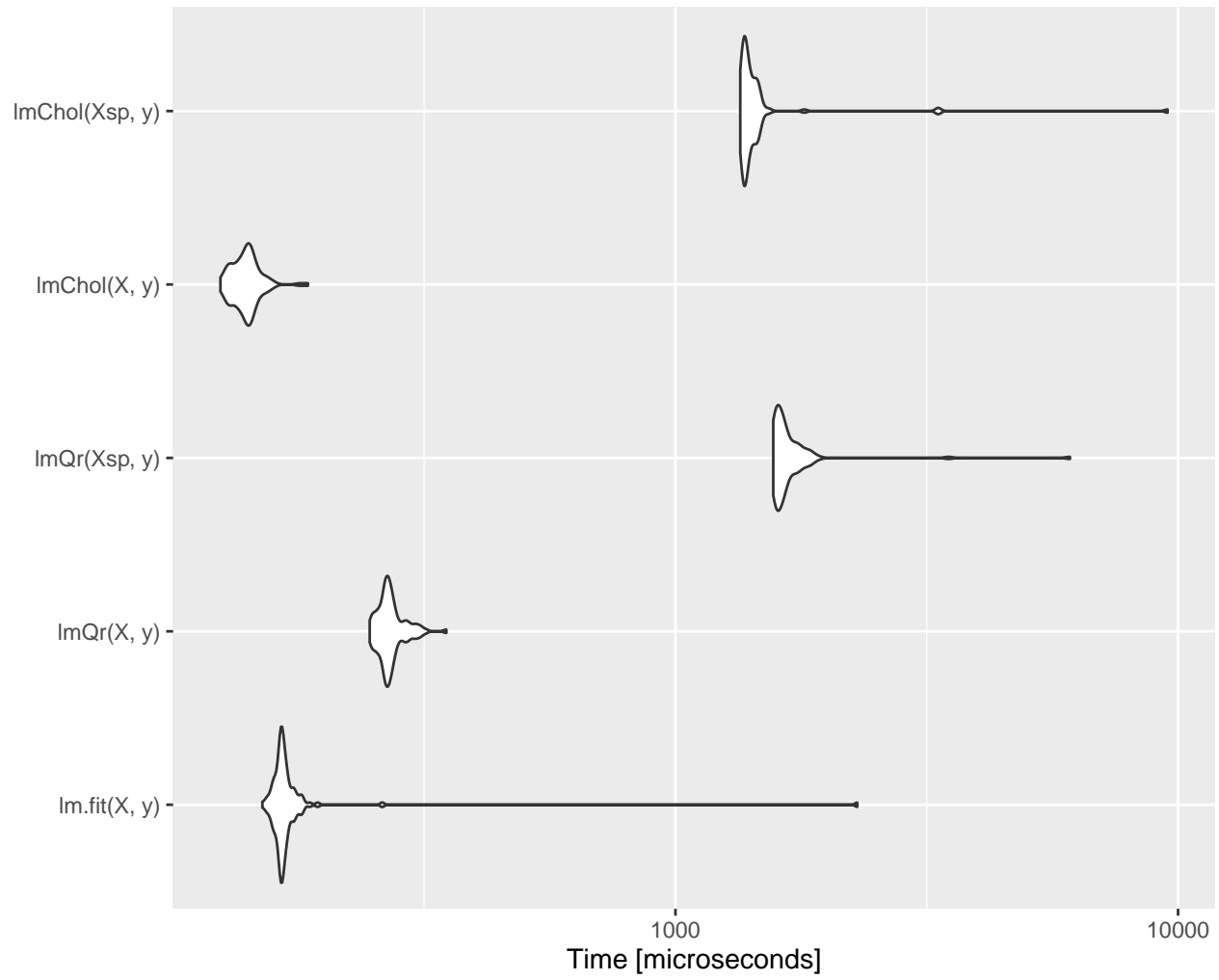


Figure 2: Performance comparison when design matrix is a sparse matrix of small dimension

1.2.3 Dense matrix of large dimension

The response Y is set as a vector with length 500. The design matrix is a 5,000 by 200 matrix. Similarly, the performance comparison is shown in Figure 3

```
set.seed(1216)
n <- 5e3
p <- 2e2
X <- matrix(rnorm(n * p), n)
y <- rowSums(X) + rnorm(n)
ggplot2::autoplot(microbenchmark(lm.fit(X, y), lmQr(X, y),
                                lmChol(X, y), times = 100))
```

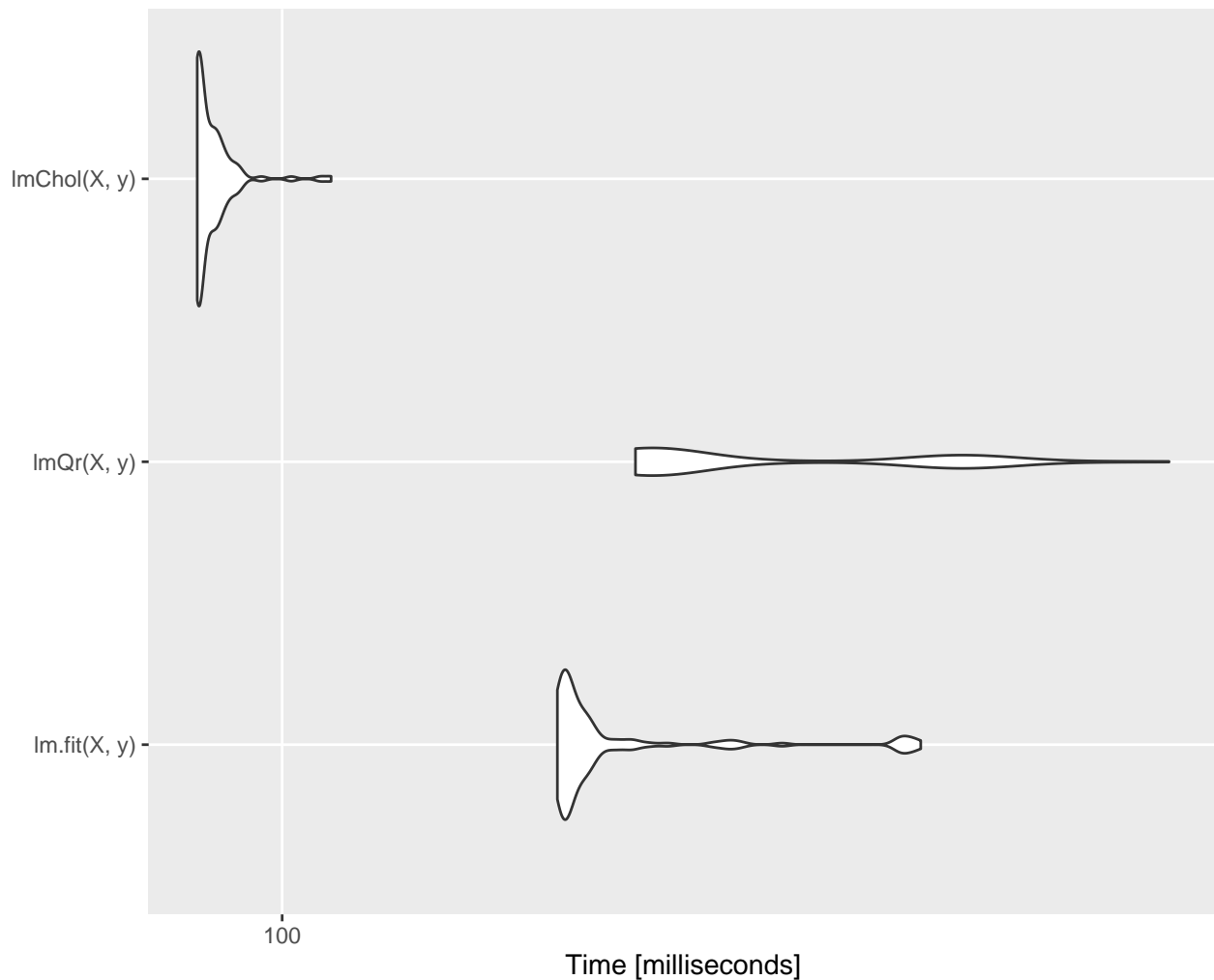


Figure 3: Performance comparison when design matrix is a dense matrix of large dimension

When a dense design matrix of large dimension is given, `lmChol` is still the fastest one. `lm.fit` is a little bit slower than `lmChol` but faster than `lmQr`.

1.2.4 Sparse matrix of large dimension

The sparsity is set to be 0.9.

```

r <- 0.9
X[sample(n * p, floor(n * p * r))] <- 0
Xsp <- methods::as(X, "dgCMatrix")
ggplot2::autoplot(microbenchmark(lm.fit(X, y), lmQr(X, y), lmQr(Xsp, y),
                                lmChol(X, y), lmChol(Xsp, y), times = 100))

```

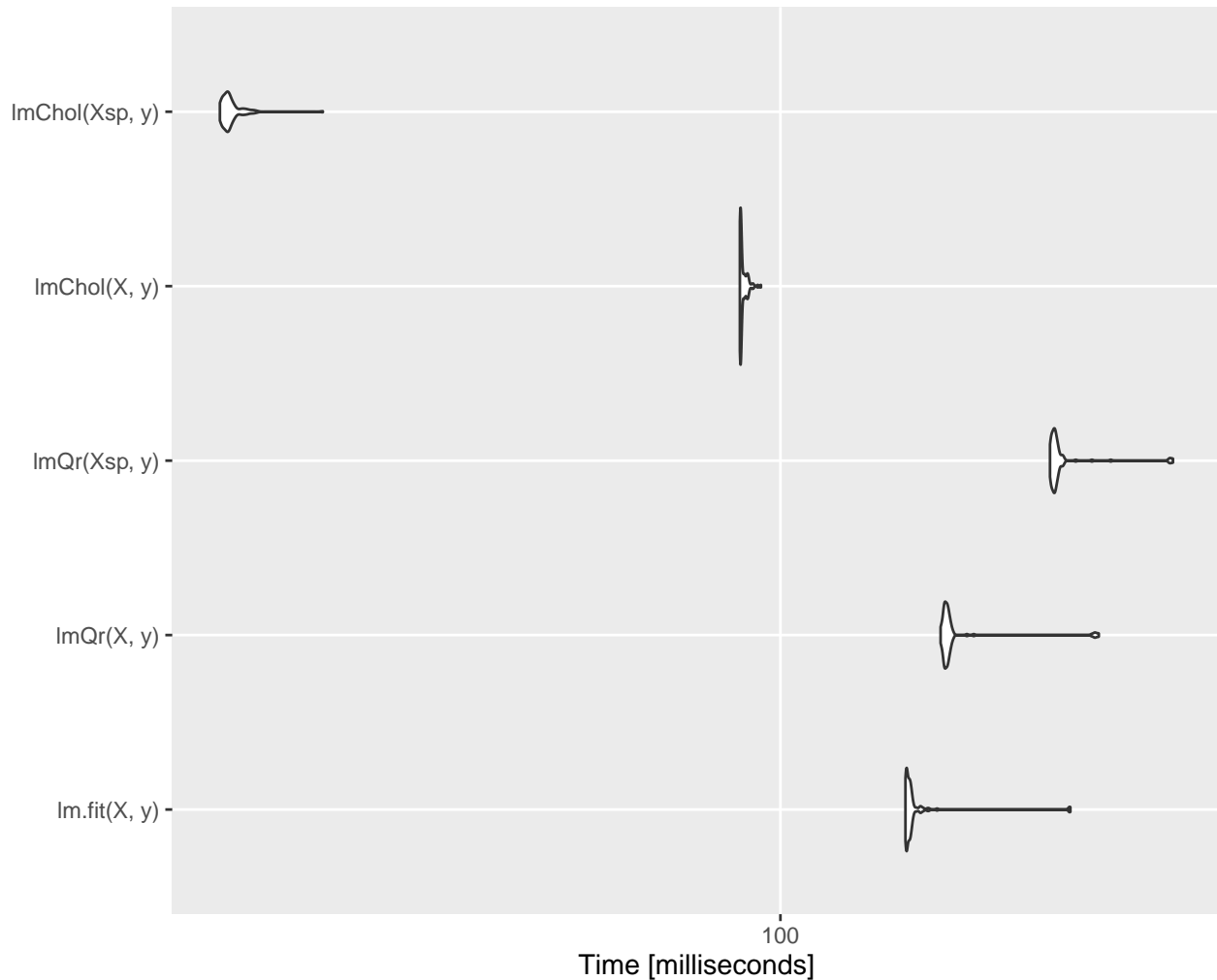


Figure 4: Performance comparison when design matrix is a sparse matrix of large dimension

From the comparison shown in Figure 4, we find that `lmChol` on sparse matrix objects becomes the fastest one. `lmChol` becomes much slower when a dense matrix object is specified, which means that the sparse matrix class object and its methods take advantage of the sparsity of the matrix and the advantage grows larger when the dimension grows. What's more, `lm.fit` is slower than `lmChol` but faster than `lmQr`. On the contrary, `lmQr` runs slower when a sparse matrix class object is input, which means those existing methods for `dgCMatrix` fail to take advantage of the sparsity after QR decomposition. (Currently, there is no `backsolve` or similar function in package **Matrix**. Package **SparseM** provides method `backsolve` for different class objects.)