

# Efficient Solution of Large Fixed Effects Problems Using R

Tom Balmat\*, Jerome P. Reiter†

January 26, 2018

## Abstract

The standard method for solving Ordinary Least Squares (OLS) regression problems in R is to use `lm()` and, for small problems, involving fewer than 1,000,000 observations and 25 independent variables, this is efficient. But for larger problems, involving tens of millions of observations and thousands of independent variables, execution of `lm()` becomes computationally impractical due to system memory requirements and execution time. Large fixed effects regression problems, involving effects with thousands of levels, present a special performance opportunity because of the large proportion of entries in the expanded design matrix (fixed effect levels translated from single columns into dichotomous indicator columns, one for each level) that are zero. For many problems, the proportion of expanded design matrix entries that are zero is above 0.995, which would be considered sparse. In this paper, we explore various methods for solving large, sparse fixed effects OLS problems and propose a method, using a combination of available R packages and custom algorithms, that efficiently computes parameter estimates and standard errors without creation of the complete expanded design matrix and limiting computations to those involving non-zero level effects. A feature, often desired in social science applications, is to estimate parameter standard errors clustered about a key identifier, such as employee ID, and large problems, with ID counts in the millions, present significant computational challenges. We present a sparse matrix Indexing algorithm that produces clustered standard error estimates that, for large fixed effects problems, is many times more efficient than standard sandwich matrix operations.

**Keywords:** fixed effects least squares solution, efficient high dimension computation, sparse matrix methods, parallel computing, clustered standard error estimation

The authors would like to thank Alex Bolton, Andres (Felipe) Barrientos, and John de Figueiredo for sharing data and models from their research and to the IT staff at Duke University's Social Science Research Institute for high performance computing support.

---

\*Social Science Research Institute, Duke University, Durham, NC 27708 (thomas.balmat@duke.edu)

†Department of Statistical Science, Duke University, Durham, NC 27708 (jerry@stat.duke.edu)

# 1 Introduction

The classic method of estimating Ordinary Least Squares (OLS) model parameters is to solve  $\mathbf{X}'\mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}'\mathbf{Y}$  for  $\hat{\boldsymbol{\beta}}$ . For small designs, where  $\mathbf{X}(n \times p)$  involves  $n$  and  $p$  on the order of  $10^4$  and  $10^1$ , respectively, this system is efficiently solved in R using the `lm()` function, while larger designs, where  $n \times p$  is on the order of  $10^7 \times 10^3$  and above, present special storage and computational challenges.<sup>1</sup> For example, creation of a  $50,000,000 \times 2,000$  numeric design matrix in R, which encodes in double floating point format (eight bytes per cell), requires  $10^{11} \times 8$  bytes, or nearly one terabyte, of memory. Further, to generate  $\mathbf{X}'\mathbf{X}$  with standard R functions, that is using `t(X)%*%X`, requires an additional copy of  $\mathbf{X}$  for the transpose, for a total memory requirement of approximately two terabytes (note that standard R functions require objects to be completely contained in on-line memory).<sup>2</sup> Physical memory constraints are easily overcome by using R x64 (the 64 bit version with increased memory addressing) installed on a system with on-line memory sufficient for the problem to be solved.<sup>3</sup> In addition to memory constraints, solution feasibility is practically limited by the processing time required to compose  $\mathbf{X}'\mathbf{X}$  and  $\mathbf{X}'\mathbf{Y}$ , solve  $\mathbf{X}'\mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}'\mathbf{Y}$  for parameter estimates, and compute  $\hat{\boldsymbol{\sigma}}^2(\mathbf{X}'\mathbf{X})^{-1}$  for (homoskedastic) standard errors. While solving large ( $10^7 \times 10^3$  and above) fixed effects regression problems, and searching for optimal methods among those currently available, it was observed that certain solutions concentrate on reduced memory usage, others target efficiency in composing  $\mathbf{X}'\mathbf{X}$  and computing parameter estimates, yet none were identified that could be considered optimal in computing standard errors. A strategy was heuristically developed that utilizes available R packages where optimal, combined with customized algorithms to improve computational efficiency where needed. This strategy is presented along with an evaluation of the efficiency of various alternatives. Additionally, an efficient algorithm is presented for estimating robust and clustered standard errors.<sup>4</sup> Traditional clustered standard error estimation involves the variance equation  $\text{Var}(\hat{\boldsymbol{\beta}}) = \text{diag}[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{U}\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}]$ , where  $\mathbf{U}$  is the  $n \times n$  matrix of within cluster residual covariances, and is estimated by setting all inter-cluster errors,  $(\boldsymbol{\epsilon}'\boldsymbol{\epsilon})_{\text{cluster}}$ , to 0 (on the assumption of inter-cluster independence).<sup>5</sup> For observation counts,  $n$ , on the order of  $10^7$  and fixed effects on the order of  $10^3$  the  $\text{Var}(\hat{\boldsymbol{\beta}})$  equation involves operations on the order of  $10^{19}$ , which is generally impracticable using standard linear algebra operations. The algorithm presented uses a sparse indexing method along with piece-wise summing of cluster row-column products in solving

<sup>1</sup>`lm()` is the base linear regression function of R (R Foundation for Statistical Computing, 2017, R Reference, `lm()`). Inspection of the program source for `lm()`, using `edit(lm())`, reveals a sequence of function calls to the R function `lm.fit()`, the C function `.Call(Cdqr1s)`, and finally the Fortran functions `dqr1s()` and `dqrdc()`. `dqr1s()` and `dqrdc()` reveal that, instead of solving  $\mathbf{X}'\mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}'\mathbf{Y}$  directly, the QR decomposition of  $\mathbf{X}$  is computed, then  $\mathbf{R}\hat{\boldsymbol{\beta}} = \mathbf{Q}'\mathbf{Y}$  is solved. For program listings see (R-core, 2017b, `Cdqr1s`) and (R-core, 2017a, `dqr1s.f`, `dqrdc.f`). Due to storage and computation inefficiencies, computing the QR decomposition of  $\mathbf{X}$  is generally impractical for problems of the size considered in this paper. QR decomposition is compared to alternative methods in *Solving the OLS Normal Equations*, below.

<sup>2</sup>The additional copy can be confirmed by monitoring total memory usage while executing the R command `sum(t(X)%*%X)`.

<sup>3</sup>Sparse matrix methods offer additional conservation by encoding non-zero matrix elements only. For more on sparse methods, see Koenker and Ng (2003) and Bates and Maechler (2017).

<sup>4</sup>Robust and clustered (heteroskedastic) standard errors are popular in social science models, where influential independent variables are correlated within groups or subjects, but not explicitly included (generally not available) in the model (for more, see Cameron and Miller (2015) and King (2015)).

<sup>5</sup>There exist methods of estimating robust standard errors that avoid direct evaluation of the above matrix equation; a popular one is known as the "Huber Sandwich Estimator" (Freedman, 2006; Zeileis, 2006; Esarey and Menger, 2017).

$\text{Var}(\hat{\beta})$  to achieve very reasonable solution times (minutes) for problems on the order of those mentioned above.

## 2 Background

While conducting research into U.S. federal employee pay disparity, the Human Capital Project at Duke University developed a requirement for fitting high dimension fixed effects OLS models to a large set of historical federal employee human capital data.<sup>6</sup> Models varied, but researchers were interested in estimates for parameters and their homoskedastic and/or heteroskedastic (clustered) standard errors. A simultaneous requirement arose from the Synthetic Data Project, also at Duke, to certify a synthetic U.S. federal human capital data set to be used in a public use verification server system.<sup>7</sup>

### 2.1 Example Data Set and Model

The U.S. Office of Personnel Management (OPM) maintains employment and human capital records for over one million non-DOD U.S. federal employees in what it calls the central personnel data file (CPDF).<sup>8</sup> CPDF “Status” records for fiscal years 1988 through 2011 were supplied by OPM to the Human Capital Project at Duke University in response to a Freedom of Information Act request.<sup>9</sup> Over 28,000,000 records were provided, each describing the human capital profile of active federal employees as of September 30 of the corresponding fiscal year; there exists one record per employee per year. Important data elements are: sex, race, age, education level, bureau employed in, occupation, and fiscal year of observation (terminology is from OPMs Guide to Data Standards).<sup>10</sup> Current research using these data have identified systematic changes in federal employee pay over time, disparities in federal employee pay by gender, and association of elections, political ideology to turnover among federal employees (Bolton and de Figueiredo, 2017; Bolton et al., 2016). For illustration, we will use a fixed effects model, similar to those used by Bolton and de Figueiredo, that measures disparity in pay by employee declared gender and race.<sup>11</sup> The model is

$$y = \beta_0 + \beta_{sex} + \beta_{race} + \beta_{sex,race} \times sex \times race + \beta_{age} \times age + \beta_{age^2} \times age^2 + \beta_{ed} \times ed_{years} + \beta_{bureau} + \beta_{occ} + \beta_{year} \quad (1)$$

where  $y$  is the logarithm of *basic pay* (OPM variable for pay) and the  $\beta$ ’s are linear effects of an employees *sex*, *race*, *sex* and *race* interaction, *age*, square of *age*, *education* (years beyond HS), *bureau* (agency employed in), *occupation*, and *fiscal year*. *Sex*, *race*, *bureau*, *occupation*, and *year* are discrete fixed effects while *age* and *education* are continuous. Fixed effects parameter estimates measure the difference in expected value of

<sup>6</sup>For more on the Human Capital Project at Duke, see Duke University Human Capital Project.

<sup>7</sup>For more on synthetic data and verification systems, see Reiter et al. (2009); Reiter (2003); Barrientos et al. (2017).

<sup>8</sup>For a general description of OPM data resources, see U.S. Office of Personnel Management (a).

<sup>9</sup>For more on CPDF Status records, see U.S. Office of Personnel Management (b).

<sup>10</sup>See OPM Guide to Data Standards, U.S. Office of Personnel Management (c).

<sup>11</sup>Bolton and de Figueiredo interact *sex* and *race* by fitting separate models for each sex.

$\log(\text{basic pay})$  between a given level and a reference level. Table 1 lists reference levels for each effect. For this model, observations are limited to full time employees with non-zero *basic pay* values and valid, non-empty values in each independent variable, giving a total included record count of 24,574,480. In constructing the design matrix  $\mathbf{X}$ , each fixed effect  $C_j$  is expanded into  $p_j - 1$  indicator columns, one for each non-reference level of  $C_j$ , where  $p_j$  is the number of distinct levels of  $C_j$ . Element (row)  $i$  of indicator column  $X_{j,k}$  corresponding to level  $k$  of fixed effect  $C_j$  contains a 1 if observation  $i$  is encoded with the  $k$ th level of  $C_j$  and contains a 0 otherwise. Table 1 lists the number of non-reference levels by effect. Including columns for the constant, continuous, and interaction variables gives a design matrix of dimension  $24,574,480 \times 1,236$ . Levels within a fixed effect are mutually exclusive since an observation can be encoded for a single level only, so that each row of the design matrix  $\mathbf{X}$  has at most one column for each fixed effect coded as a 1, with the remainder necessarily containing 0. Large  $p_j$  values cause overall low density (proportion of 1 entries) in  $\mathbf{X}$  and it is this property that sparse methods exploit in efficient construction of  $\mathbf{X}'\mathbf{X}$  and solution of the OLS normal equations. Table 1 also lists mean density by fixed effect (average density of columns associated with each effect). Overall design matrix density in the OPM data set relative to model (1) is 0.006, making it quite sparse and a good candidate for algorithmic and solution comparison.

Table 1: Pay Disparity Fixed Effect Reference Levels, Number of Non-ref Levels, and Density

Effect	Reference Level	Non-Ref Levels	Density
Sex	M (Male)*	1	0.4800
Race	E (White)*	4	0.0821
Bureau	114009000 (Veterans Administration)**	409	0.0020
Occupation	303 (Miscellaneous Clerk and Assistant)**	791	0.0012
Year	1988***	23	0.0420

\* customary reference levels for pay equity research

\*\* highest marginal frequency

\*\*\* first year in study data

### 3 Design Matrix and Operation Density

To solve  $\mathbf{X}'\mathbf{X}\hat{\beta} = \mathbf{X}'\mathbf{Y}$ , we need  $\mathbf{X}'\mathbf{X}$ . Suppose that we are given a large, sparse design matrix  $\mathbf{X}$  ( $25,000,000 \times 2,000$  with a high proportion of elements equal to 0) along with a list of elements (row  $i$ , column  $j$ ) that are non-zero. An efficient alternative to generating  $\mathbf{X}'\mathbf{X}$  with matrix operations is to accumulate products of pairs of non-zero elements (indicated by  $i$  and  $j$ ) into corresponding positions of the resulting  $\mathbf{X}'\mathbf{X}$  matrix.<sup>12</sup> Since  $\mathbf{X}'\mathbf{X}$  is symmetric, additional efficiency is gained by accumulating the upper triangle only, then copying its transpose to the lower triangle. Sparse methods imply limiting computations strictly

<sup>12</sup>The standard method in R is `Z <- t(X)%*%X`.

to those that affect the final result, in the case of matrix multiplication, products that are non-zero. The *Matrix and Operation Density* appendix establishes that matrix density is strongly influenced by fixed effect dimension, that operation density is a result of matrix density, and that theoretical measures of operation density indicate an opportunity, for sufficiently sparse problems, to eliminate a high proportion of standard matrix operations in producing  $\mathbf{X}'\mathbf{X}$ .<sup>13</sup> Figure 2 (*Operation Density*) of the appendix indicates that for models involving fixed effects with approximately 50 or more levels, the proportion of effective operations (those that must be executed) approaches zero. Models with fixed effects in this range are very common (50 or more occupations, states, counties, schools, etc.), making implementation of a solution method both relevant and appealing. Model (1), combined with the example OPM data set, has an operation density of approximately 0.00001, making it a candidate for computational efficiency improvement.<sup>14</sup>

## 4 Review of Available Solutions

Many statistical software programs recognize fixed effects, or categorical variables, in solving regression problems. Popular solutions such as `proc glm` from SAS (SAS Institute, 2017, `proc GLM`), the `regress` command of Stata (Stata Corporation, 2017, `regress`), and the `lm()` function of R (R Foundation for Statistical Computing, 2017, `lm()`) implement a common method of expanding categorical variables into columns of the design matrix, one for each level, eliminating the column corresponding to a user declared or system selected reference level to avoid linear dependence, and finally solving the normal equations  $\mathbf{X}'\mathbf{X}\hat{\beta} = \mathbf{X}'\mathbf{Y}$  for  $\hat{\beta}$ .<sup>15</sup> The *Review of Available Solutions* appendix evaluates popular software solutions for solving large OLS regression problems.<sup>16</sup> Since this paper targets an R audience, solutions are limited to R functions and packages available at time of writing, including `lm()`,<sup>17</sup> `biglm()`,<sup>18</sup> `bigglm()`,<sup>18</sup> `biglm.big.matrix()`,<sup>19</sup> `SparseM.slm.fit()`,<sup>20</sup> `speedlm()`,<sup>21</sup> `lfe()`,<sup>22</sup> and combined functions from the `Matrix` and `Matrix Models` packages.<sup>23</sup> Each addresses efficiency obstacles in different ways, some targeting memory constraints only, others addressing both efficient computation and use of memory. Some employ parallel processing methods, others create external operating system files when a design matrix exceeds on-line memory capacity. Each solves or approximates solution of the OLS normal equations, providing estimates for all continuous and fixed effect parameters. Some provide homoskedastic standard error estimates, some provide matrix operations for

---

<sup>13</sup>The *Matrix and Operation Density* appendix is available in the on-line repository (Duke University Synthetic Data Project, Fixed effects solution git repository).

<sup>14</sup>Operation density can be derived from a constructed  $\mathbf{X}'\mathbf{X}$  matrix using the fixed effect indicator product entries, since they report the number of intersecting non-zero elements and indicate the number of pseudo-multiplication operations performed

<sup>15</sup>Verification of fixed effect expansion into binary columns with that of the reference level omitted can be accomplished with the `model.matrix()` function in R (R Foundation for Statistical Computing, 2017, R Reference, `model.matrix()`) and the `base` option of the `regress` command in Stata (Stata Corporation, 2017, `regress`).

<sup>16</sup>Available in the on-line repository (Duke University Synthetic Data Project, Fixed effects solution git repository).

<sup>17</sup>See (R Foundation for Statistical Computing, 2017, R Reference, `lm()`).

<sup>18</sup>See Lumley (2015).

<sup>19</sup>See Emerson and Kane (2016).

<sup>20</sup>See Koenker and Ng (2003).

<sup>21</sup>See Enea et al. (2017).

<sup>22</sup>See Gaure (2016).

<sup>23</sup>These include `sparse.model.matrix()`, `lm.fit.sparse()`, and `solve()` (Bates and Maechler, 2017, 2015).

estimating standard errors using  $\hat{\sigma}^2(\mathbf{X}'\mathbf{X})^{-1}$ , while few offer estimates of heteroskedastic robust or clustered standard errors. Example data sets and models are employed to illustrate performance features of various solutions. Additionally, model (1) is fit, when possible, using the example OPM data set described above.

While testing available solutions for feasibility and efficiency, with a requirement to meet project deadlines, the authors developed algorithms and programs in R and C for efficiently indexing  $\mathbf{X}$ , composing  $\mathbf{X}'\mathbf{X}$ , solving for parameter estimates, and computing homoskedastic, robust, and clustered standard errors.<sup>24</sup> The custom R function `feXTX()` implements this solution and performance of it is compared to other methods in *Performance of Solutions with Simulated Data and Models* and *Performance with OPM Data and Models*, below.

An alternative to solving the OLS equations,  $\mathbf{X}'\mathbf{X}\hat{\beta} = \mathbf{X}'\mathbf{Y}$ , for all  $\hat{\beta}$  is to re-specify the problem, transforming continuous independent variables to deviations from within-level means. This is referred to as *demeaned regression*. Although efficient, the method does not provide estimates for fixed effect parameters or standard errors for any parameter. It is mentioned for completeness, covering the case where only continuous variable parameter estimates are desired.<sup>25</sup>

## 5 Solving the OLS Normal Equations

It will be seen in later sections on performance that, for large fixed effects problems, the most efficient methods involve solution of  $\mathbf{X}'\mathbf{X}\hat{\beta} = \mathbf{X}'\mathbf{Y}$  for parameter estimates,  $\hat{\beta}$ , and computation of  $\hat{\sigma}^2(\mathbf{X}'\mathbf{X})^{-1}$  for standard errors (as opposed to some form of approximation). We see that, in addition to  $\mathbf{X}'\mathbf{X}$ ,  $\mathbf{X}'\mathbf{Y}$  is needed but, as demonstrated in the *Efficient Indexing* and *Review of Available Solutions* appendices, it is efficiently produced with sparse indexing or matrix multiplication operations.

Of the various methods for solving  $\mathbf{X}'\mathbf{X}\hat{\beta} = \mathbf{X}'\mathbf{Y}$  for  $\hat{\beta}$  two that are well established and known for computational efficiency are QR decomposition and Cholesky decomposition.<sup>26</sup> <sup>27</sup> The efficiency of QR and Cholesky decomposition derive from a strategy of forming pairs of matrices, with at least one triangular, such their product represents the left hand matrix in the system being solved,  $\mathbf{X}'\mathbf{X}$  in our case.<sup>28</sup> In the case of Cholesky decomposition, an upper triangular matrix  $\mathbf{R}$  is computed such that  $\mathbf{R}'\mathbf{R} = \mathbf{X}'\mathbf{X}$ , then  $\mathbf{R}'\mathbf{R}\hat{\beta} = \mathbf{X}'\mathbf{Y}$  is solved in two stages:  $\mathbf{R}'\mathbf{W} = \mathbf{X}'\mathbf{Y}$  for  $\mathbf{W}$  then  $\mathbf{R}\hat{\beta} = \mathbf{W}$  for  $\hat{\beta}$ . The triangular nature of  $\mathbf{R}$

<sup>24</sup>These are presented in appendices *Efficient Indexing of X in Composing X'X*, *An Efficient X'X Indexing Algorithm in R*, *Parallel Cholesky Decomposition Algorithm*, and *X<sup>-1</sup> Using Parallel Cholesky Decomposition*, available in the on-line repository (Duke University Synthetic Data Project, Fixed effects solution git repository).

<sup>25</sup>For more on demeaned regression, see Gormley and Matsa (2013) and Williams (2015).

<sup>26</sup>For more on QR decomposition, see Wikipedia (2017b) and Vandenberghe (2017). QR decomposition is implemented by the R function `qr()` (R Foundation for Statistical Computing, 2017, `qr()`) and is the default algorithm used by `lm()`.

<sup>27</sup>For more on Cholesky decomposition, see Wikipedia (2017a) and Heath (2013). Cholesky decomposition is implemented by the R function `chol()` (R Foundation for Statistical Computing, 2017, `chol()`).

<sup>28</sup>An alternate method, as used by `lm()`, is to compute the QR decomposition of  $\mathbf{X}$ , such that  $\mathbf{QR} = \mathbf{X}$ , where  $\mathbf{Q}$  is  $(n \times p)$  orthogonal and  $\mathbf{R}$  is  $(p \times p)$  triangular, giving  $\mathbf{X}'\mathbf{X}\hat{\beta} = \mathbf{X}'\mathbf{Y} \implies \mathbf{R}'\mathbf{Q}'\mathbf{Q}\mathbf{R}\hat{\beta} = \mathbf{R}'\mathbf{Q}'\mathbf{Y} \implies \mathbf{R}'\mathbf{R}\hat{\beta} = \mathbf{R}'\mathbf{Q}'\mathbf{Y} \implies \mathbf{R}\hat{\beta} = \mathbf{Q}'\mathbf{Y}$ , assuming  $\mathbf{R}$  is invertible. However, for even small fixed effects problems  $(1,000,000 \times 500)$ , the computational cost of computing  $\mathbf{Q}$  and  $\mathbf{R}$  is many times that of composing  $\mathbf{X}'\mathbf{X}$  then computing its QR decomposition to solve  $\mathbf{X}'\mathbf{X}\hat{\beta} = \mathbf{X}'\mathbf{Y}$ . For large problems, as mentioned in the introduction, storage of  $\mathbf{Q}$   $(n \times p)$  becomes prohibitive. For a comparison of the efficiency of computing the QR decomposition of  $\mathbf{X}$  vs. that of  $\mathbf{X}'\mathbf{X}$ , see the *QR Decomposition of X vs. X'X* appendix, available in the on-line repository (Duke University Synthetic Data Project, Fixed effects solution git repository).

makes solution of associated equations very efficient by enabling the use of computationally simple forward and backward solving steps. In addition to facilitating solution of  $\mathbf{X}'\mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{R}'\mathbf{R}\hat{\boldsymbol{\beta}} = \mathbf{X}'\mathbf{Y}$ ,  $\mathbf{R}$  can be used to solve  $\mathbf{X}'\mathbf{X}\mathbf{W} = \mathbf{I}$ , where  $\mathbf{W} = (\mathbf{X}'\mathbf{X})^{-1}$ , the computationally significant component of the homoskedastic parameter covariance equation  $\text{Cov}(\hat{\boldsymbol{\beta}}) = \sigma^2 (\mathbf{X}'\mathbf{X})^{-1}$ . QR decomposition follows a similar strategy, but with matrices  $\mathbf{Q}$  and  $\mathbf{R}$ , where  $\mathbf{QR} = \mathbf{X}'\mathbf{X}$ ,  $\mathbf{Q}$  is orthogonal, and  $\mathbf{R}$  is upper triangular. In practice, the computational effort required to generate  $\mathbf{R}$  (or  $\mathbf{Q}$  and  $\mathbf{R}$ ) combined with solving two simpler systems tends to be less than in solving  $\mathbf{X}'\mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}'\mathbf{Y}$  using direct Gaussian elimination and the resulting equations and solution tend to exhibit improved numerical stability.<sup>29</sup> Although both QR decomposition and Cholesky decomposition produce near mathematically equivalent results for problems of the size considered, it has been the authors' experience that, for large fixed effect regression problems, Cholesky decomposition requires less time to solve both systems of equations, one for  $\hat{\boldsymbol{\beta}}$  and one for  $\text{Cov}(\hat{\boldsymbol{\beta}})$ , than does QR decomposition. As an example, table 2 compares times to compute the QR decomposition, Cholesky decomposition, and inverse of  $\mathbf{X}'\mathbf{X}$  from two design matrices (one of size 5,000,000  $\times$  6,200; the other 10,000,000  $\times$  16,000) using the standard R functions `qr()`, `qr.solve()`, `chol()`, and `chol2inv()`.<sup>30 31 32</sup> The apparent four to one performance ratio gives Cholesky decomposition a significant advantage.

Table 2: Time (in minutes) to Compute Decomposition and Inverse of  $\mathbf{X}'\mathbf{X}$

Method	6,200 level matrix	16,000 level matrix
<code>qr() + qr.solve()*</code>	10.0	172.7
<code>chol() + chol2inv()</code>	2.5	44.1

\* computed using the slightly more efficient LAPACK option, as opposed to the default LINPACK

We will restrict our discussion to Cholesky decomposition. While solving various large fixed effects regression problems, it was observed that as much as three-fourths of the entire computation cycle [composition of  $\mathbf{X}'\mathbf{X}$ , computation of  $\mathbf{R}$ , solution of  $\hat{\boldsymbol{\beta}}$ , solution of  $\text{Cov}(\hat{\boldsymbol{\beta}})$ ] was devoted to the final three items, all involving  $\mathbf{R}$ . Although, algorithmically,  $\mathbf{R}$  can be computed in simultaneous parallel operations, the base R functions `chol()` and `chol2inv()` are implemented as a sequence of row and column operations, with a single row or column being computed before subsequent rows or columns are begun.<sup>33</sup> As part of the current project, two Cholesky related functions were developed: one to compute the decomposition  $\mathbf{R}$ , given a composed  $\mathbf{X}'\mathbf{X}$ , and one to compute  $(\mathbf{X}'\mathbf{X})^{-1}$  using  $\mathbf{R}$ . Both are implemented in C, using the R package `Rcpp`<sup>34</sup>, and employ parallel methods, using `OpenMP`<sup>35</sup>. Source listings are given in appendices *Parallel Cholesky Decomposition*

<sup>29</sup>See Higham (1990).

<sup>30</sup>Data sets were generated as described in *Performance of Solutions with Simulated Data and Models*, below.

<sup>31</sup>`qr.solve(X)` computes the inverse of  $\mathbf{X}$  (R Foundation for Statistical Computing, 2017, `qr.solve()`).

<sup>32</sup>`chol2inv(R)` computes the inverse of  $\mathbf{X} = \mathbf{R}'\mathbf{R}$  (R Foundation for Statistical Computing, 2017, `chol2inv()`).

<sup>33</sup>This is verified by observing processor activity during execution.

<sup>34</sup>See Eddelbuettel et al. (2017).

<sup>35</sup>Available in `Rtools` (Ripley and Murdoch, 2017).

*Algorithm and Parallel Algorithm to Compute Inverse of  $\mathbf{X}$  Using the Cholesky Decomposition.*<sup>36</sup> Function `choleskyDecomp()` computes  $\mathbf{R}$  and function `cholInvDiag()` computes  $(\mathbf{X}'\mathbf{X})^{-1}$ . Figure 1 compares the performance of the custom parallel algorithms to that of their base R counterparts, solid lines indicate execution times of custom functions, dashed lines for standard R functions.<sup>37</sup> With performance of approximately eight times that of base R functions, equating to a savings of over five hours for the largest problem tested, the benefit of custom, parallel algorithm implementation is significant.

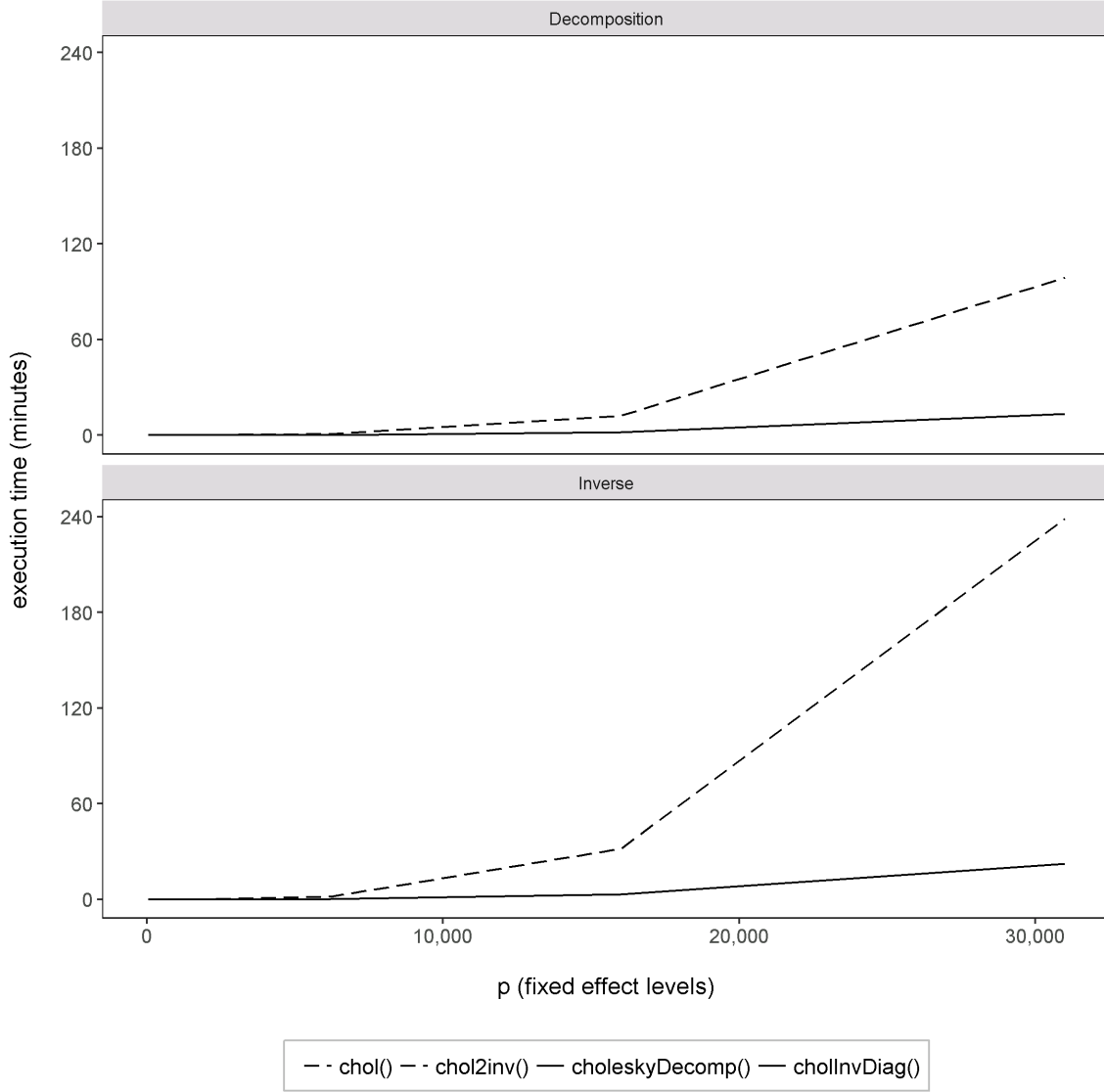


Figure 1: Comparison of execution times to compute Cholesky decomposition and  $(\mathbf{X}'\mathbf{X})^{-1}$  using custom parallel algorithm and base R functions. Five independent data sets were generated for each level of  $p_k$  (number of fixed effect levels). Observation counts ranged from 10,000 (small  $p_k$ ) to 25,000,000 (large  $p_k$ ).

<sup>36</sup> Available in the om-line repository (Duke University Synthetic Data Project, Fixed effects solution git repository).

<sup>37</sup> Sample data sets were generated as described in *Performance of Solutions with Simulated Data and Models*, below.



To verify consistency of results between the custom parallel and base R functions, element-wise differences in computed  $\mathbf{R}$  or  $(\mathbf{X}'\mathbf{X})^{-1}$  matrices of absolute value greater than  $10^{-10}$  were tested for using `which(abs( $\mathbf{R}_1 - \mathbf{R}_2$ ) > 1e-10)` and `which(abs(( $\mathbf{X}'\mathbf{X}$ )1-1 - ( $\mathbf{X}'\mathbf{X}$ )2-1) > 1e-10)`, where  $\mathbf{R}_1$  and  $\mathbf{R}_2$  are custom and base R Cholesky decomposition results, respectively, and  $(\mathbf{X}'\mathbf{X})_1^{-1}$  and  $(\mathbf{X}'\mathbf{X})_2^{-1}$  are custom and base R inverses, respectively. None were reported. *Note that the  $1e^{-10}$  threshold is subjective in that, two values on the order of  $1e^{-12}$  could have a difference less than  $1e^{-10}$  but, based on their magnitude, could be quite different.*

## 6 Performance of Solutions with Simulated Data and Models

We now compare the performance and results of methods presented in the *Review of Available Solutions* section. Before using the example OPM data set, since it is not available to the public, and to give the reader a ready means of testing methods presented, simulated data sets with specifiable covariate relationships will be used. The *Simulated Data Generation and Tests* appendix contains an R script that produces a data frame, `feDat`, with one dependent (continuous numeric) vector  $Y$  and independent vectors  $X_1$ ,  $X_2$ ,  $X_3$ ,  $X_4$ , and  $X_5$ , where  $X_1$  and  $X_2$  are continuous numeric vectors and  $X_3$ ,  $X_4$ , and  $X_5$  are categorical (character) fixed effects.<sup>38</sup>  $Y$  is simulated as a linear function of the  $X$  values plus normally distributed error such that

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \beta_{3i} + \beta_{4i} + \beta_{5i} + \epsilon \quad (2)$$

where  $\beta_{3i}$ ,  $\beta_{4i}$ , and  $\beta_{5i}$  are the effects corresponding to levels of  $X_3$ ,  $X_4$ , and  $X_5$  of observation  $i$ . Values for the  $\beta$ 's and  $\epsilon$  are available in the appendix. Table 3 lists test parameters (number of observations and number of levels for  $X_3$ ,  $X_4$ , and  $X_5$ ) used to generate simulated data sets for evaluation of execution times.

Also in the *Simulated Data* appendix are functions for fitting model (2) to simulated data, using `lm()`, `biglm()`, `bigglm()`, `biglm.big.matrix()`, `biglm.big.matrix()`, `SparseM.slm.fit()`, `speedlm()`, `lfe()`, `feXTX()`, `Matrix1`, and `Matrix2`.<sup>39 40 41</sup> Note that appropriate packages must be installed in order to execute related model fitting functions. Figure 2 plots mean execution time by method using five independently simulated data sets for parameter sets 1 through 9. In addition to execution efficiency, the total on-line memory used by an algorithm is important, since it is generally limited and shared by multiple concurrent users and processes. As a quick assessment of memory required by each method, table 4 lists maximum differential of memory in use (from baseline prior to execution, as reported by Windows task

<sup>38</sup>The appendix is available in the on-line repository (Duke University Synthetic Data Project, Fixed effects solution git repository).

<sup>39</sup>*Fitting* consists of generating coefficient and (homoskedastic) standard error estimates for all continuous and fixed effect independent variables.

<sup>40</sup>`feXTX()` employs efficient indexing and the custom parallel Cholesky algorithms `choleskyDecomp()` and `cholInvDiag()`, as earlier presented, along with base R functions `forwardsolve()` and `backsolve()`.

<sup>41</sup>`Matrix1` employs `sparse.model.matrix()`, `lm.fit.sparse(method="cholesky")`, and `solve()` from the `Matrix` and `Matrix Models` packages. `Matrix2` employs `sparse.model.matrix()` from the `Matrix Models` package, base R functions `forwardsolve()` and `backsolve()`, and custom parallel Cholesky algorithms `choleskyDecomp()` and `cholInvDiag()`.

Table 3: Simulated Data Set Parameters for Execution Time Evaluation

Parameter Set	n (Observations)	Levels $X_3$	Levels $X_4$	Levels $X_5$
1	10,000	10	15	20
2	25,000	10	20	40
3	50,000	15	30	50
4	100,000	15	40	60
5	150,000	25	40	60
6	250,000	25	50	75
7	500,000	50	75	100
8	1,000,000	75	150	250
9	1,500,000	100	250	500
10	2,000,000	100	500	1,000
11	5,000,000	200	1,000	5,000
12	10,000,000	1,000	5,000	10,000
13	25,000,000	1,000	5,000	25,000

manager) while the algorithms were executed using a single instance of data generated by parameter set 9 (the highest dimension set that all methods were executed with).

Table 4: Maximum On-line Memory Used While Fitting Simulation Parameter Set 9

Method	Memory Used (Gb)
<code>lm()</code>	19.3
<code>biglm()</code>	19.3
<code>bigglm()</code>	19.4
<code>biglm.big.matrix()</code>	0.80*
<code>SparseM.slm.fit()</code>	20.3
<code>speedlm()</code>	28.8
<code>lfe()</code>	10.4
<code>feXTX()</code>	1.0
<code>Matrix<sub>1</sub></code>	1.0
<code>Matrix<sub>2</sub></code>	1.0

\* backing file created

Efficient execution is important, but accuracy of results is of the utmost importance. Considering the computational challenges of fitting high dimension models to large data sets, obtaining an authoritative set of parameter estimates for objective comparison is problematic. As an alternative, we might take the

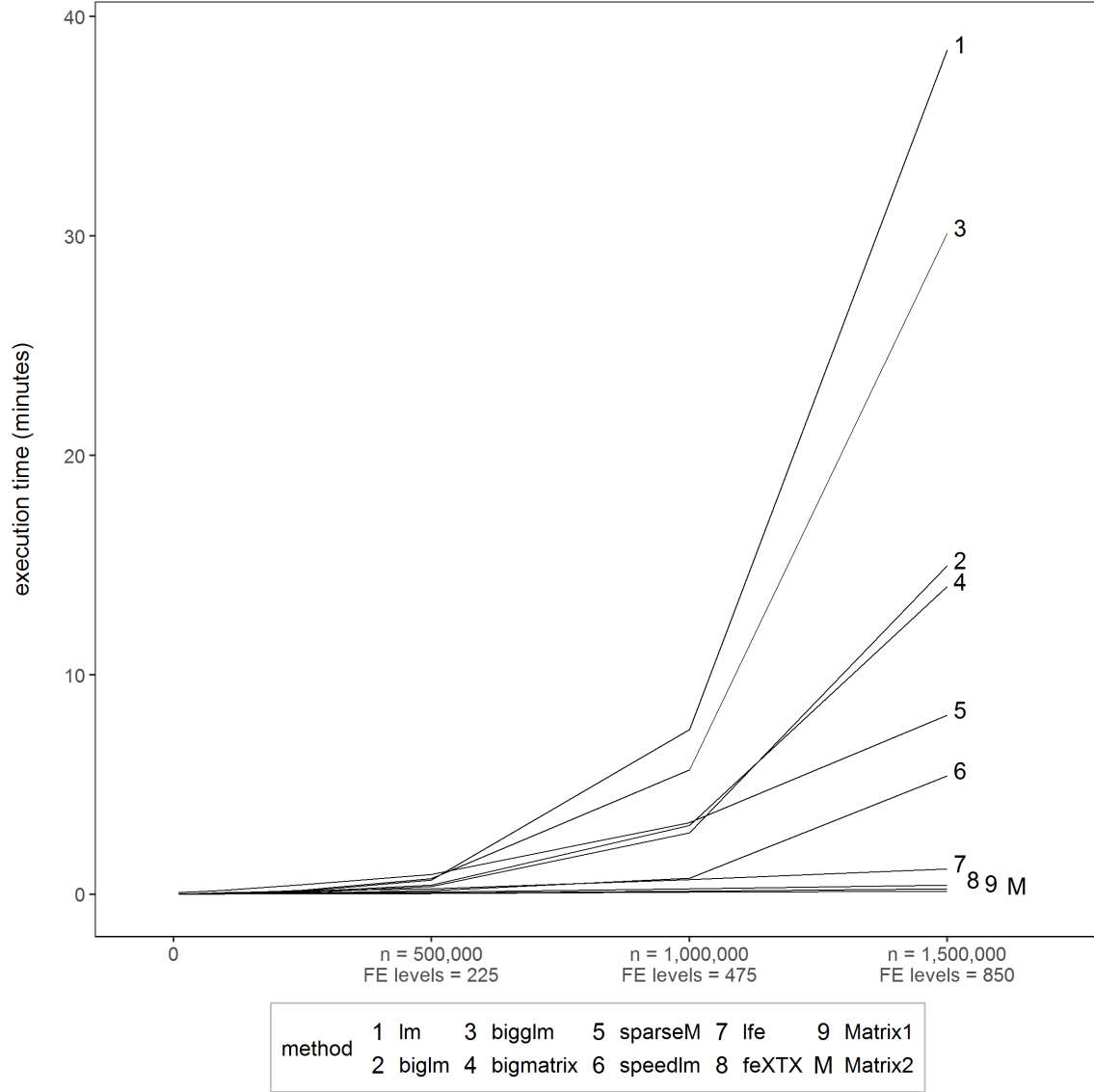


Figure 2: Mean execution time to solve graduated fixed effects problems by method. Execution time is the differential of “elapsed” time as reported by `proc.time()`. All processing was performed on a single dedicated server.

approach of comparing estimates to a consensus of those from established methods. If no major deviations are observed then we can state that, with data sets and models tested, our method introduces no error beyond that observed in methods currently in use. Table 5 shows, by method, the number of parameter estimates, using simulation parameter set 9, that differ with those generated by `feXTX()`, grouped by absolute value ranges listed in the column headers.<sup>42</sup> There does seem to be consensus, with all methods except `lfe()`

<sup>42</sup>The smallest magnitude of all `feXTX()` parameter estimates, for this simulated data set, was 0.5. Absolute deviations from this value of less than  $10^{-7}$  should be considered insignificant. Comparison functions are available in the *Simulated Data* appendix, available in the on-line repository (Duke University Synthetic Data Project, Fixed effects solution git repository).

having all deviations of parameter estimates less than  $10^{-7}$ .<sup>43</sup> <sup>44</sup>

Table 5: Comparison of Parameter Estimates, Evaluated Solutions vs. `feXTX()`, Simulated Data Set 9

Method	$0 < 10^{-12}$	$10^{-12} < 10^{-10}$	$10^{-10} < 10^{-7}$	$10^{-7} < 10^{-5}$	$> 10^{-5}$
<code>lm()</code>	384	464	2	0	0
<code>biglm()</code>	341	507	2	0	0
<code>bigglm()</code>	341	507	2	0	0
<code>bigmatrix()</code>	341	507	2	0	0
<code>sparseM()</code>	0	101	749	0	0
<code>speedlm()</code>	0	101	749	0	0
<code>lfe()</code>	0	0	3	0	847
<code>Matrix<sub>1</sub></code>	0	101	749	0	0
<code>Matrix<sub>2</sub></code>	0	101	749	0	0

In terms of execution and memory efficiency, `lfe()`, `feXTX()`, `Matrix1`, and `Matrix2` appear to significantly outperform the other solutions and maintain efficiency throughout the range of smaller data sets. Continuing with more demanding data sets (parameter sets 10, 11, 12, and 13) and limiting our study to `lfe()`, `feXTX()`, `Matrix1`, and `Matrix2`, we see in figure 3 a continued divergence of execution time, with `feXTX()` and `Matrix2` exhibiting near linear increase with respect to problem size, executing the largest problem, with 25,000,000 observations and 31,000 fixed effect levels, in one-twentieth to one-half the time required by `lfe()` and `Matrix1`. Memory requirements to fit model (2) using parameter set 13 were ??? Gb for `lfe()`, ??? Gb for `feXTX()`, ??? Gb for `Matrix1`, and ??? Gb for `Matrix2`.

## 7 Performance with OPM Data and Models

However compelling the case for efficiency is made using simulated data, ultimately we must assess utility with actual data and models. Performance and issues regarding execution of `lm()`, `biglm()`, `bigglm()`, `biglm.big.matrix()`, `SparseM.slm.fit()`, and `speedlm()` using the OPM data and models are discussed in the *Review of Available Solutions* appendix.<sup>45</sup> Due to obstacles (errors during execution) and lengthy execution times with larger models, these solutions are not considered here. As previously stated, the example OPM data set contains 24,574,480 observations and two moderate dimension fixed effects (bureau

<sup>43</sup>Although `feXTX()` and `Matrix2` use identical Cholesky and forward/back solve functions, the difference in estimates is traced to a difference in  $X'X$  construction, `feXTX()` using indicator columns for counts and sums of indexed columns, `Matrix2` executing traditional matrix multiplication.

<sup>44</sup>By default, `lfe()` does not estimate an intercept and fixed effect reference levels are not specifiable when the faster “kaczmarz” method is used (the method appears to force highest frequency levels as reference). This causes differences in remaining estimates. Differences in effects between non-reference level and desired reference level could be computed from supplied estimates, but this was not pursued. The alternative “cg” method permits reference level specification but, in trials conducted with data and models presented, is many times slower than “kaczmarz.” For more on this, see (Gaure, 2016, pg. 26).

<sup>45</sup>Available in the on-line repository (Duke University Synthetic Data Project, Fixed effects solution git repository).

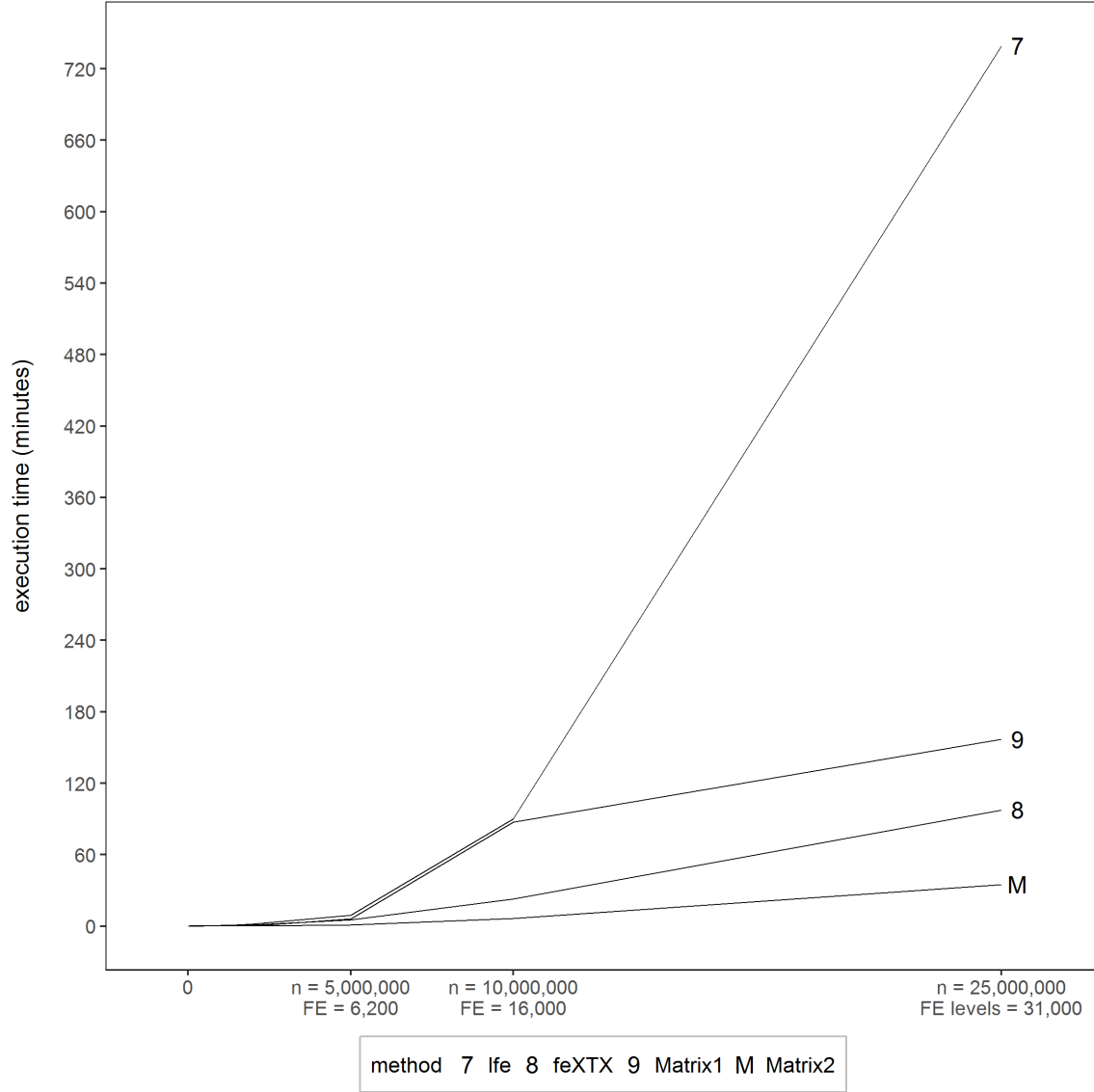


Figure 3: Mean execution time to solve graduated fixed effects problems, including high dimension effects. Limited to `lfe()`, `feXTX()`, `Matrix1`, and `Matrix2` methods. Execution time is the differential of “elapsed” time as reported by `proc.time()`. All processing was performed on a single dedicated server.

and occupation, with 409 and 791 levels, respectively). A fully expanded design matrix would have dimension  $24,574,480 \times 1,236$  and relatively low density (matrix density of 0.0006, operation density of 0.00001). Problem size, density, and the fact that actual area research was delayed due to difficulties obtaining parameter and standard error estimates to hypothesized models in reasonable time, make the data set and model (1) good candidates for analysis. Execution times and memory requirements to fit model (1) and compute

homoskedastic standard errors, using the example OPM data, are listed in table 6.<sup>46</sup> <sup>47</sup>

Table 6: Execution Time and Memory Required to Fit Model (1) to OPM Data

Method	Execution Time (min)	Memory used (Gb)
<b>lfe()</b>	34.8	128*
<b>feXTX()</b>	10.6	12
<b>Matrix<sub>1</sub></b>	6.0	8
<b>Matrix<sub>2</sub></b>	3.8	8

\* exceeds memory available on 64 Gb test server; execution accomplished on alternate dedicated server with 256 Gb of memory; “non-estimable function” warnings observed on alternate server

As with the simulated data sets and models, the **Matrix<sub>2</sub>** method significantly outperforms other methods evaluated, with no penalty in memory usage. As previously mentioned, **Matrix<sub>2</sub>** consists of sparse matrix functions (importantly `sparse.model.matrix()`) from the **Matrix** and **Matrix Models** packages combined with custom Cholesky decomposition and inverse functions. While searching for efficient solutions to large computational problems, it is important to observe and measure performance of individual components and to develop methods for combining optimal elements into an overall optimal solution. Table 7 lists, by method, absolute deviations of coefficient estimates to those computed by **feXTX()**. With a minimum estimate of 0.0003 (absolute value, computed by **feXTX()**), deviations of less than  $10^{-7}$  indicate at least three significant digits of precision, making deviations substantively insignificant.<sup>48</sup> As previously mentioned, restrictions on reference level specification cause significant deviations in certain estimates computed by **lfe()**.<sup>49</sup>

Table 7: Comparison of Parameter Estimates, Evaluated Solutions vs. **feXTX()**, OPM Data

Method	$0 < 10^{-12}$	$10^{-12} < 10^{-10}$	$10^{-10} < 10^{-7}$	$10^{-7} < 10^{-5}$	$> 10^{-5}$
<b>lfe()</b>	2	145	1,022	27	36
<b>Matrix<sub>1</sub></b>	0	1	1,231	0	0
<b>Matrix<sub>2</sub></b>	0	1	1,231	0	0

<sup>46</sup>Although the OPM data are not available for public release, the appendix *Measuring Performance of Fixed Effect Solutions Using OPM Data*, available in the on-line repository (Duke University Synthetic Data Project, Fixed effects solution git repository), contains the R instructions and function calls used to generated this table.

<sup>47</sup>For **feXTX()** execution, *sex*  $\times$  *race* interaction columns were constructed and declared as continuous variables, since it has been observed that, for interactions of low dimension fixed effects (*sex* has two levels, *race* has five), this approach is generally more efficient than having the indexing algorithm compose interaction columns. Often, with experimentation, observation, and measurement, algorithm features or alternate model specification can be employed for efficiency gains.

<sup>48</sup>This has been confirmed in discussions with members of the Human Capital project at Duke.

<sup>49</sup>On inspection, it is observed that **lfe()** chooses highest frequency levels as reference levels. By coincidence, the Human Capital models also use highest frequency levels for Bureau and Occupation. This enables **lfe()** to produce estimates that are nearer to those of **feXTX()**, **Matrix<sub>1</sub>**, and **Matrix<sub>2</sub>** with models fit to the OPM data than with models fit to simulated data.

$$y = \beta_0 + \beta_{sex} + \beta_{race} + \beta_{sex,race} \times sex \times race + \beta_{age} \times age + \beta_{age} \times age^2 + \beta_{ed} \times ed_{years} + \beta_{bureau} + \beta_{occ} + \beta_{year} + \beta_{sta} \quad (3)$$

Table 8: Execution Time and Memory Required to Fit Model (3) to OPM Data

Method	Execution Time (min)	Memory used (Gb)
<b>feXTX()</b>	10.6	20
<b>Matrix<sub>1</sub></b>	6.0	12
<b>Matrix<sub>2</sub></b>	3.8	12

## 8 Robust and Clustered Standard Errors

To compensate for heteroskedasticity of within-group errors, analysts often report parameter estimates with respect to robust or clustered standard errors (Cameron and Miller, 2015). **feXTX()** computes robust and clustered standard errors using the analytical equation for  $\text{Cov}(\hat{\beta})$  derived from the OLS normal equations:

$$\text{Cov}(\hat{\beta}) = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{u}\mathbf{u}'\mathbf{X} (\mathbf{X}'\mathbf{X})^{-1} \quad (4)$$

where  $\mathbf{u}\mathbf{u}'$  is the  $n \times n$  matrix of within cluster residual covariances, and is estimated by setting all inter-cluster errors,  $(\epsilon'\epsilon)_{cluster}$ , to 0 (on the assumption of inter-cluster independence).<sup>50</sup> Note that  $\mathbf{u}\mathbf{u}'$  is  $n \times n$  which, for the problems we consider, can be quite large,  $24,574,480 \times 24,574,480$  in our example OPM data set. Further, it is involved in operations with another large matrix,  $\mathbf{X}(24,574,480 \times p)$ . Standard matrix operations involving objects of this size are computationally impractical, in both time and memory requirements. However, sparse indexing methods similar to those used by **feXTX()** to compose  $\mathbf{X}'\mathbf{X}$  can be used to compute  $\text{Cov}(\hat{\beta})$  very efficiently.<sup>51</sup> In fact, having constructed sparse indices to fixed effect columns of  $\mathbf{X}$  along with  $(\mathbf{X}'\mathbf{X})^{-1}$ , we simply proceed with construction of  $\mathbf{U}$  (which is efficiently accomplished in parallel), then multiplication of matrices in equation (4) referencing only non-zero elements of  $\mathbf{X}$  (which is also done in parallel). Note that, since the standard errors of  $\hat{\beta}$  are generally of interest, only  $\text{Var}(\hat{\beta}) = \text{diag} [\text{Cov}(\hat{\beta})]$  is computed.

*Robust standard error algorithm outline:*

<sup>50</sup>For more on heteroskedasticity and inter-cluster independence, see Esarey and Menger (2017).

<sup>51</sup>These are presented in appendices *Efficient Indexing of X in Composing X'X* and *An Efficient X'X Indexing Algorithm in R*, available in the on-line repository (Duke University Synthetic Data Project, Fixed effects solution git repository).

- Compute robust, uncorrelated, heteroskedastic (independent, non-identically-distributed) errors
- Estimate robust variances, square of  $SE(\hat{\beta})$ , as  $\text{Var}(\hat{\beta}) = \text{diag}[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{u}\mathbf{u}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}]$ , where  $\mathbf{u}$  is the vector of observation response to predicted value errors,  $\mathbf{u} = \mathbf{Y} - \hat{\mathbf{Y}}$
- This derives from the expression for parameter variances from the normal OLS equations
- $\mathbf{u}\mathbf{u}'$  forms an  $n \times n$  diagonal  $\epsilon$  variance-covariance matrix with all off-diagonal elements set to 0, which asserts the assumption of independent (uncorrelated) errors
- But, since the diagonal elements are expected to be non-constant, this method models parameter standard errors in a heteroskedastic (independent, but not identically distributed) error setting
- Since  $\mathbf{X}'\mathbf{u}\mathbf{u}'\mathbf{X}$  is symmetric, construct upper triangle only then copy transpose to lower triangle

*Clustered standard error algorithm outline:*

- Compute heteroskedastic, correlated within cluster, independent (uncorrelated) between cluster, non-identically-distributed (heteroskedastic) errors
- Estimate clustered standard errors using  $\text{Var}(\hat{\beta}) = \text{diag}[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{u}\mathbf{u}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}]$ , where  $\mathbf{u}\mathbf{u}'$  is the var-cov matrix of observation to estimate errors,  $\mathbf{u}\mathbf{u}' = \text{Cov}(\mathbf{Y} - \hat{\mathbf{Y}})$
- This is the expression for  $\text{Var}(\hat{\beta})$  from ordinary least squares estimates of beta on the assumption that  $\mathbf{X}$  is non-stochastic and  $\mathbf{u}\mathbf{u}'$  is an estimate of error covariance and, further, that covariance observed within clusters (groups) estimates that in the population while covariance between groups is 0
- $\mathbf{u}\mathbf{u}'$  is altered such that all off-diagonal inter-group elements, that is  $\mathbf{u}\mathbf{u}'_{ij}$  where indices  $i$  and  $j$  belong to different groups, are set to 0
- The modified (inter-group independent error)  $\mathbf{u}\mathbf{u}'$  has the form:

$$\begin{bmatrix} \mathbf{u}\mathbf{u}'_{group_1} & 0 & 0 & \cdots & 0 \\ 0 & \mathbf{u}\mathbf{u}'_{group_2} & 0 & \cdots & 0 \\ 0 & 0 & \mathbf{u}\mathbf{u}'_{group_3} & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & \mathbf{u}\mathbf{u}'_{group_k} \end{bmatrix}$$

where  $\mathbf{u}\mathbf{u}'_{group_i}$  is the  $n_{group_i} \times n_{group_i}$  sub-matrix of the  $\epsilon$  covariance matrix corresponding to group  $i$

- Note that with this construction of  $\mathbf{u}\mathbf{u}'$ ,  $\mathbf{X}'\mathbf{u}\mathbf{u}'\mathbf{X}$  is the sum (over all groups  $i$ ) of  $\mathbf{X}'_{group_i} \mathbf{u}\mathbf{u}'_{group_i} \mathbf{X}_{group_i}$  where  $\mathbf{X}_{group_i}$  and  $\mathbf{u}\mathbf{u}'_{group_i}$  are the rows of  $\mathbf{X}$  and sub-matrix of  $\mathbf{u}\mathbf{u}'$  corresponding to group  $i$  [to see this, imagine multiplying  $\mathbf{X}'$  by  $\mathbf{u}\mathbf{u}'$  and consider the effect of  $\mathbf{u}\mathbf{u}'$  elements of rows that do not



correspond to a particular group - each resulting column corresponds to products from a single group - then multiplying that by  $\mathbf{X}$  (each row belongs to a single group) gives a  $p \times p$  sum of individual group products]

- Note that  $\mathbf{X}'_{group_i} \mathbf{u}_{group_i} = [\mathbf{u}'_{group_i} \mathbf{X}_{group_i}]'$ , making  $\mathbf{X}'_{group_i} \mathbf{u} \mathbf{u}'_{group_i} \mathbf{X}_{group_i}$  symmetric
- Also,  $\mathbf{u}'_{group_i} \mathbf{X}_{group_i}$  is a  $1 \times p$  vector, where  $p$  is the column dimension of  $\mathbf{X}$  (the design matrix)
- For each group ID, construct  $\mathbf{v} = \mathbf{u}'_{group_i} \mathbf{X}_{group_i}$  then use it to compute  $\mathbf{v}' \mathbf{v} = \mathbf{X}'_{group_i} \mathbf{u} \mathbf{u}'_{group_i} \mathbf{X}_{group_i}$

Because `lfe()` requires specification of estimable functions to compute estimates for an intercept and non-reference level fixed effects, which is not pursued, comparison of robust and clustered standard errors is made to those computed by Stata, which is known for computational efficiency and accuracy. Table 9 lists execution times to fit model (1) to the OPM data and to compute robust and clustered standard errors using `lfe()`, `feXTX()`, and Stata (Stata/MP version 13.1). `lfe()` is included for execution time comparison. All processing was executed on a single, dedicated, 24 core Windows 7 server. The approximate 2 to 1 performance ratio of `feXTX()` to both `lfe()` and Stata indicates an advantage to computing heteroskedastic robust and clustered standard errors by solving the analytical variance equation, eq (4), using efficient indexing methods. The largest (absolute value) pair-wise deviation in Stata and `feXTX()` computed standard errors was less than  $10^{-5}$ , lending credence to accurate analytical equation evaluation.

Table 9: Execution Time (in minutes) to Fit Model (1) to OPM Data and Compute Robust and Clustered Standard Errors

Method	Robust SE Time (min)	Clustered SE Time (min)	Memory used (Gb)
<code>lfe()</code> *	—	51.9	128
<code>feXTX()</code> **	13.2	13.2	20/38
Stata***	51.5	55.5	<20

\* bootstrap estimates computed; robust standard errors require use of ancillary bootstrap estimation functions which was not pursued; executed on alternate server with additional required memory  
message received: "`chol.default() - matrix is either rank-deficient or indefinite`"  
message received: "`warning; non-estimable function`"

\*\* SE's result from evaluation of analytical standard error equation

\*\*\* Using Stata/MP Version 13.1

## 9 Further Development

While developing `feXTX()` and associated algorithms, several opportunities for improvement of performance or utility were identified but, due to practical limitations of time and striving to hold a course of completing primary features, they are withheld for future development. Among them are:

- Improved sub-setting and transmission of fixed effect vectors to parallel cores transmit only levels to be operated on by a given core
- Elimination of fixed effect level index vector export (parallel R, under Windows; an Rcpp/OpenMP implementation would enable shared memory)
- Parallelization of interaction vector-products
- Implement a hybrid solution where  $\mathbf{X}'\mathbf{X}$  is constructed and normal equations are solved using sparse methods from the **Matrix** package, followed by homoskedastic, robust, and clustered standard error estimation using custom parallel Cholesky and heteroskedastic SE methods demonstrated above
- Implement a parallel algorithm for forward and back solving of OLS system of equations using computed Cholesky decomposition

## 10 Conclusion

As large data sets become increasingly available and available data become increasingly complex, new methods are being developed to model systems and identify patterns using both deep and broad statistical probing. At the same time, traditional modeling methods, with tried and true histories, remain relevant and in use, and with data sets from the “small data” period of their origin, familiar methods remain computationally practical. However, with larger data sets, computation of analytical equations from traditional models may be of such order to make solution impractical without introducing estimation or convergence methods. In this paper, we have demonstrated that the traditional, analytical normal equations of ordinary least squares (OLS) in a high dimension fixed effects setting can be effectively solved by employing a combination of efficient indexing of sparse vectors, parallel computation, and targeted C programming. A methodology was presented that 1.) identifies sources of computational inefficiency in standard OLS functions as implemented in R; 2.) evaluates the efficiency and accuracy of existing computational solutions; 3.) develops alternative strategies and algorithms to efficiently compute results that are mathematically equivalent to those produced by standard functions; and 4.) adapts solutions beyond basic estimates of model parameters to computationally intensive estimates of parameter variance. Efficiency methods and algorithms targeting large data sets and models presented in the paper include:

- Analysis and measurement of design matrix and operation density
- Alternative representation of sparse fixed effect indicator columns using compact index vectors
- Construction of  $\mathbf{X}'\mathbf{X}$  using efficient alternatives to vector multiplication
- Solution of parameter estimates from the OLS normal equations ( $\mathbf{X}'\mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}'\mathbf{Y}$ ) using a custom, parallelized Cholesky decomposition algorithm implemented in C

- Solution of homoskedastic parameter standard errors from the analytical equation  $\text{Cov}(\hat{\beta}) = \hat{\sigma}^2(\mathbf{X}'\mathbf{X})^{-1}$ , using a custom, parallelized function, implemented in C, that computes the inverse of a matrix from its Cholesky decomposition
- Solution of heteroskedastic, robust and clustered, parameter standard errors from the analytical equation  $\text{Cov}(\hat{\beta}) = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{u}\mathbf{u}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}$ , using sparse fixed effect index columns produced during construction of  $\mathbf{X}'\mathbf{X}$

Finally, algorithm efficiency was compared to that of eight solutions and packages available for R using small to large simulated data sets. Performance of the most efficient solutions was evaluated using a moderate sized data set and models taken from actual research in U.S federal pay practices. In all cases, the custom algorithms presented out-performed the remaining methods, eliminating hours to days of compute time. Accuracy of results was confirmed using a consensus of results from all solutions and, in some cases, by making comparisons to estimates produced by alternative software (Stata). Although restricted to a limited family of problems (fixed effects OLS), it is hoped that the general approach of inefficiency isolation, experimentation, algorithm development, performance comparison, and result certification demonstrated here may be useful to readers when facing computationally challenging problems encountered in their research.

## References

- Barrientos, A. F., Bolton, A., Balmat, T., Reiter, J. P., de Figueiredo, J. M., Machanavajjhala, A., Chen, Y., Kneifel, C., and DeLong, M. A Framework for Sharing Confidential Research Data, Applied to Investigating Differential Pay by Race in the U. S. Government. Tech. rep., National Bureau of Economic Research Working Paper 23534, 2017.
- Bates, D. and Maechler, M. R Matrix Models Package, 2015. URL <https://cran.r-project.org/web/packages/MatrixModels/MatrixModels.pdf>.
- Bates, D. and Maechler, M. R Matrix Package, 2017. URL <https://cran.r-project.org/web/packages/Matrix/Matrix.pdf>.
- Bolton, A. and de Figueiredo, J. M. Measuring and explaining the gender wage gap in the U.S. federal government. Tech. rep., Duke University Law School, 2017.
- Bolton, A., de Figueiredo, J. M., and Lewis, D. E. Elections, ideology, and turnover in the U.S. federal government. Tech. rep., National Bureau of Economic Research Working Paper 22932, 2016.
- Cameron, A. C. and Miller, D. L. A Practitioners Guide to Cluster-Robust Inference. *Journal of Human Resources*, 50(2):317–332, 2015.
- Duke University Human Capital Project. URL <https://ssri.duke.edu/projects/human-capital-career-dynamics-and-organizational-performance-us-federal>.
- Duke University Synthetic Data Project. Fixed effects solution git repository. URL <https://github.com/DukeSynthProj/FixedEffectsSolutionResources>.
- Eddelbuettel, D., Francois, R., Allaire, J., Ushey, K., Kou, Q., Russell, N., Bates, D., and Chambers, J. R Rcpp Package, 2017. URL <https://cran.r-project.org/web/packages/Rcpp/Rcpp.pdf>.
- Emerson, J. W. and Kane, M. J. R biganalytics Package, 2016. URL <https://cran.r-project.org/web/packages/biganalytics/biganalytics.pdf>.
- Enea, M., Meiri, R., , and Kalimi, T. R speedglm Package, 2017. URL <https://cran.r-project.org/web/packages/speedglm/speedglm.pdf>.
- Esarey, J. and Menger, A. Practical and Effective Approaches to Dealing with Clustered Data, 2017. URL <http://jee3.web.rice.edu/cluster-paper.pdf>.
- Freedman, D. A. On the so-called "huber sandwich estimator" and "robust standard errors". *The American Statistician*, 60(4):299–302, 2006.
- Gaure, S. R lfe Package, 2016. URL <https://cran.r-project.org/web/packages/lfe/lfe.pdf>.
- Gormley, T. A. and Matsa, D. A. Common Errors: How to (and Not to) Control for Unobserved Heterogeneity, 2013. URL <http://portal.idc.ac.il/en/main/research/caesareacenter/annualsummit/documents/gormley-matsa.pdf>.
- Heath, M. T. Parallel Numerical Algorithms, 2013. URL [https://courses.engr.illinois.edu/cs554/fa2013/notes/07\\_cholesky.pdf](https://courses.engr.illinois.edu/cs554/fa2013/notes/07_cholesky.pdf).
- Higham, N. J. Analysis of the Cholesky Decomposition of a Semi-definite Matrix. 1990. URL [http://eprints.ma.man.ac.uk/1193/01/covered/MIMS\\_ep2008\\_56.pdf](http://eprints.ma.man.ac.uk/1193/01/covered/MIMS_ep2008_56.pdf).
- King, G. How Robust Standard Errors Expose Methodological Problems They Do Not Fix, and What to Do About It. *Political Analysis*, 23:159–179, 2015.
- Koenker, R. and Ng, P. SparseM: A Sparse Matrix Package for R, 2003. URL <http://www.econ.uiuc.edu/~roger/research/sparse/SparseM.pdf>.

- Lumley, T. R biglm Package, 2015. URL <https://cran.r-project.org/web/packages/biglm/biglm.pdf>.
- R-core. R Parallel Package, 2016. URL <https://stat.ethz.ch/R-manual/R-devel/library/parallel/doc/parallel.pdf>.
- R-core. R Project Source, 2017a. URL <https://svn.r-project.org/R/trunk/src/appl>.
- R-core. R Project Source, lm.c function, 2017b. URL <https://svn.r-project.org/R/trunk/src/library/stats/src/lm.c>.
- R Foundation for Statistical Computing. R Reference, 2017. URL <https://cran.r-project.org/doc/manuals/r-release/fullrefman.pdf>.
- Reiter, J. P. Model diagnostics for remote access regression servers. *Statistics and Computing*, 13:371–380, 2003.
- Reiter, J. P., Oganian, A., and Karr, A. F. Verification servers: Enabling analysts to assess the quality of inferences from public use data. *Computational Statistics and Data Analysis*, 53:1475–1482, 2009.
- Ripley, B. and Murdoch, D. Rtools, 2017. URL <https://cran.r-project.org/bin/windows/Rtools/>.
- SAS Institute. Proc GLM, 2017. URL <https://support.sas.com/documentation/cdl/en/statug/63033/HTML/default/viewer.htm#glm.toc.htm>.
- Stata Corporation. regress Command, 2017. URL <http://www.stata.com/manuals13/rregress.pdf>.
- U.S. Office of Personnel Management. Data, Analysis, and Documentation, a. URL <https://www.opm.gov/policy-data-oversight/data-analysis-documentation/>.
- U.S. Office of Personnel Management. FedScope, b. URL [https://www.fedscope.opm.gov/datadefn/aehri\\_sdm.asp](https://www.fedscope.opm.gov/datadefn/aehri_sdm.asp).
- U.S. Office of Personnel Management. Guide to Data Standards, c. URL <https://catalog.data.gov/dataset/guide-to-data-standards-gds>.
- Vandenberghe, L. QR Factorization, 2017. URL <http://www.seas.ucla.edu/~vandenbe/133A/lectures/qr.pdf>.
- Wikipedia. Cholesky decomposition, 2017a. URL <https://pdfs.semanticscholar.org/f229/a57ee5611ca84a8936fd9c29a3f1f19dc1e9.pdf>.
- Wikipedia. QR Decomposition, 2017b. URL [https://en.wikipedia.org/wiki/QR\\_decomposition](https://en.wikipedia.org/wiki/QR_decomposition).
- Williams, R. Panel Data: Very Brief Overview, 2015. URL <https://www3.nd.edu/~rwilliam/stats2/Panel.pdf>.
- Zeileis, A. Object-oriented computation of sandwich estimators. *Journal of Statistical Software*, 16(9), 2006.