

Least-Angle Regression and LASSO for Large Datasets

Chris Fraley and Tim Hesterberg

Technical Report

Insightful Corporation

October 1, 2007; revised April 28, 2008

Abstract

Least-Angle Regression and the LASSO (ℓ_1 -penalized regression) offer a number of advantages in variable selection applications over procedures such as stepwise or ridge regression, including prediction accuracy, stability and interpretability. We discuss formulations of these algorithms that extend to datasets in which the number of observations could be so large that it would not be possible to access the matrix of predictors as a unit in computations. Our methods require a single pass through the data for orthogonal transformation, effectively reducing the dimension of the computations required to obtain the regression coefficients and residual sums-of-squares to the number of predictors, rather than the number of observations.

Keywords: regression, regularization, ℓ_1 penalty, lasso, scalable, massive datasets, tall datasets.

1 Introduction

This paper addresses the extension of least-angle regression (Efron et al. 2004) and the LASSO method (Tibshirani 1996) to linear models with very large numbers of observations. We write the model as:

$$y = \begin{pmatrix} \mathbf{1} & X \end{pmatrix} \begin{pmatrix} \beta \\ b \end{pmatrix} + \epsilon, \quad (1.1)$$

where y is a numeric response, $\mathbf{1}$ a column of ones, X an $n \times p$ numeric matrix of predictors, β an intercept, b a vector of p additional regression coefficients, and ϵ a vector of residuals. We assume that:

- The number of observations n may be so large that X (and possibly y) could not be stored in memory.
- The number of predictors p is sufficiently small that matrices with dimensions of order p could be held in memory and used in computations.
- X and y are stored in such a way that the rows of X can be accessed in a sequential blockwise fashion, and the corresponding components of the response y can be accessed with each block of rows of X . This scenario is typical of data storage in commercially-available databases.

In this situation, computation time is dominated by the hardware interface, hence the need for a one-pass strategy. Least-angle regression (LAR) and its extension to LASSO were originally described in terms of X and y (Efron et al. 2004), and an implementation suitable for datasets of

moderate size has been made available (Efron and Hastie 2003). Because this formulation requires direct access to X and y throughout the computation, it is not practical when the data resides out of memory.

The main idea in this paper is to do all of the necessary calculations in memory using an orthogonal transformation of the matrix of predictors and the response that reduces the dimension of the computations. The reduced matrices and vectors can be accumulated in a single pass through the data, which need not reside in memory, but is assumed accessible in blocks of rows. The dimension-reduction process is a memory-efficient version of the QR decomposition that has been used for linear regression on large datasets that are sequentially blockwise accessible. Least-angle regression and its LASSO extension involve varying sets of predictors, and we also make use of updating techniques for the QR factorization to accommodate subsets of predictors in linear regression. Finally, we show that it is possible to obtain the correlations needed for selecting the predictors, and do all of the other computations required for LAR and LASSO, after the initial factorization and transformation, without accessing the out-of-memory data.

The paper is organized as follows. The basic LARS algorithm is described in terms of X and y in Section 2. We show in Section 3 how to express all of the necessary computations without accessing X or y after initial factorization and transformation to matrices and vectors with dimensions of order p . In Section 4 we describe how to achieve this factorization and transformation using sequential blockwise access to rows of X along with simultaneous access to corresponding entries of the response y . The results of Section 3 and Section 4 allow extension of LARS to very large blockwise-accessible datasets. Section 5 summarizes the results and provides further discussion. Derivations of some of the relations used in the algorithms are given in the appendix, as well as information on a S-PLUS implementation.

2 Algorithm Description

Efron et al. (2003, 2004) develop least-angle regression as a variable-selection technique for linear models. There are three basic variants: *least-angle regression (LAR)*, *ℓ_1 penalty (LASSO)*, and *forward stagewise*, which are collectively referred to as *LARS*. Although the methods described here can be applied to all of these algorithms, we discuss only LAR and LASSO since the forward stagewise variant can require many steps and is not widely used.

In Efron et al. (2003, 2004), the predictors and response are transformed to have mean 0. This effectively assumes an intercept in the model, as shown in the appendix. Efron et al. (2004) also assume that the predictors have unit Euclidean length, but we do not make this assumption here. We write the model as

$$\tilde{y} = \tilde{X}b + \epsilon, \tag{2.2}$$

where \tilde{y} denotes the response with its mean subtracted, and \tilde{X} denotes the matrix of predictors X , each with its mean subtracted. Later we show how to implement the methods in reduced dimensions, without first centering or scaling the variables.

The methods proceed iteratively in a series of steps. We make the following definitions:

- b_k : the vector of coefficients at (the end of) step k ,
- \mathcal{A}_k : the active set of predictors during step k ,
- p_k : the number of active predictors,
- \tilde{X}_k : the columns of \tilde{X} corresponding to \mathcal{A}_k ,

- P_k : a $p_k \times p$ matrix of zeroes and ones such that $\tilde{X}_k = \tilde{X}P_k^T$,
- $r_k \equiv \tilde{y} - \tilde{X}b_k$: the vector of partial residuals at step k .

Initially all coefficients b_0 are zero, and \mathcal{A}_0 is empty. A series of models is fitted in which predictors are successively added to or dropped from the active set, and coefficients are updated. The final model has a maximal set of linearly-independent predictors.

At the k th step, the active set is updated, and a new set of coefficients b_k are obtained by determining a step length λ_k along a direction d_k from b_{k-1} :

$$b_k \leftarrow b_{k-1} + \lambda_k d_k, \quad (2.3)$$

with $0 \leq \lambda_k \leq 1$. The active set and step length calculations differ for LAR and LASSO, and we compare these below.

The direction d_k comes from the least-squares estimate based on the active set. Let

$$z_k = \underset{z}{\operatorname{argmin}} \|y - \tilde{X}_k z\| \quad (2.4)$$

be the least-squares coefficients based on the current active set. Equivalently, z_k satisfies the normal equations

$$\tilde{X}_k^T \tilde{X}_k z_k = \tilde{X}_k^T y. \quad (2.5)$$

The direction d_k is the vector of length p which is equal to $z_k - P_k b_{k-1}$ in the entries corresponding to the active predictors, and 0 elsewhere:

$$d_k = P_k^T (z_k - P_k b_{k-1}).$$

The step length λ_k to be taken from b_{k-1} along d_k to get the new set of coefficients is determined by the algorithm (LAR or LASSO), based on the correlations of the predictors with the partial residuals:

$$\operatorname{cor}(\tilde{X}, r_k) \equiv \frac{\mathcal{D}_X^{-1} \tilde{X}^T r_k}{\|r_k\|_2},$$

where

$$\mathcal{D}_X \equiv \operatorname{diag} \left(\|\tilde{X}_{*1}\|_2, \|\tilde{X}_{*2}\|_2, \dots, \|\tilde{X}_{*p}\|_2 \right).$$

The correlations are proportional to the inner products between the partial residuals and predictors, standardized by their norms. Dropping the denominator for simplicity, we work with

$$c_k = \mathcal{D}_X^{-1} \tilde{X}^T r_k \quad (2.6)$$

and with the more general

$$c(b) = \mathcal{D}_X^{-1} \tilde{X}^T (\tilde{y} - \tilde{X}b).$$

At the start of a step, the active set is defined to be the predictors corresponding to the correlations that have the largest magnitude. These correlations would vanish if a unit step were taken in the direction d_k . The first predictors to enter the active set are those with the largest absolute correlations with \tilde{y} . Subsequently, the differences in step length and active variable selection for the different LAR and LASSO are:

- **Least-Angle Regression**

The step length λ_k in LAR is the smallest step such that one or more predictors that are not in the current active set \mathcal{A}_k have correlation equal in magnitude to the correlations of members of \mathcal{A}_k at $b_{k-1} + \lambda d_k$. Those predictor(s) are added to the active set for the following step. The required computations are straightforward because

$$c(b_{k-1} + \lambda d_k) = \mathcal{D}_X^{-1} \tilde{X}^T [\tilde{y} - \tilde{X}(b_{k-1} + \lambda d_k)] \quad (2.7)$$

is a linear function of λ . In LAR, predictors never leave the active set once they are added.

- **LASSO and ℓ_1 Penalty Regression**

A *LASSO* (Tibshirani 1996) solution minimizes

$$\min_b \|\tilde{y} - \tilde{X}_k b\|_2^2 + \theta_k \|\mathcal{D}_X^{-1} b\|_1$$

for some $\theta_k > 0$. Coefficients are scaled in the ℓ_1 penalty term for consistency with Tibshirani (1996) and Efron et al. (2004), where columns of \tilde{X} are normalized.

At a LASSO solution, correlations corresponding to non-zero components of b_k are maximal in magnitude and have sign equal to the corresponding element of b_k (e.g. Osborne et al. 2000). In the LASSO modification of least-angle regression, the step length λ_k is the smallest step such that either $b_{k-1} + \lambda d_k$ changes sign for one or more of the predictors in the current active set, or the LAR step criterion holds (Efron et al. 2004). In either case, $b_k = b_{k-1} + \lambda_k d_k$ satisfies the LASSO optimality conditions. When the LASSO step is determined by a sign change, the predictors at which the coefficients change sign are zero and are dropped from the active set at the next iteration (no new predictors are added).

3 Alternative Implementation in Reduced Dimensions

In this section, we show how to implement LAR and LASSO using a Cholesky factor of $\tilde{X}^T \tilde{X}$ and the corresponding transformation of \tilde{y} , effectively reducing the dimensions of the computations. This yields the same sequence of coefficients (up to numerical error) that would have been obtained with the original formulation, but does not require access to X and y after initial factorization and transformation. Later, in Section 4, we describe the one-pass blockwise transformation of $\begin{pmatrix} \mathbf{1} & X \end{pmatrix}$ and y to the forms needed for the methods described here.

3.1 Cholesky and QR Factorizations

A $p \times p$ upper triangular matrix \tilde{R} satisfying

$$\tilde{R}^T \tilde{R} = \tilde{X}^T \tilde{X},$$

is called a Cholesky factor¹ of \tilde{X} . The Cholesky factor \tilde{R} can be obtained by applying a series of orthogonal Householder transformations to \tilde{X} . There is an associated QR factorization of \tilde{X} in which

$$\tilde{X} = \tilde{Q} \tilde{R},$$

¹The Cholesky factor is unique up to the signs of the rows.

where \tilde{Q} is an $n \times p$ matrix with orthogonal columns. In ordinary least squares, the coefficients b for the regression

$$\tilde{y} = \tilde{X}b + \epsilon$$

satisfy the normal equations

$$\tilde{X}^T \tilde{X}b = \tilde{X}^T \tilde{y},$$

or equivalently

$$\tilde{R}b = \tilde{Q}^T \tilde{y}.$$

The latter has advantages for numerical accuracy, and also for efficiency in applications like LAR and LASSO. It is not necessary to form the matrix \tilde{Q} , since the product $\tilde{Q}^T y$ can be accumulated by sequential application of the transformations used to form \tilde{R} . For extensive discussion of the QR factorization and its use in regression computations, see e.g. Chapter 3 of Demmel (1997), Chapter 2 of Trefethen and Bau (1997), or Chapter 2 of Leader (2004). In Section 4, we show describe how to obtain the Cholesky factor via sequential blockwise application of Householder transformations.

3.2 Directions and Correlations

Let \tilde{R} be an upper triangular Cholesky factor of \tilde{X} , and let

$$\check{y} \equiv \tilde{Q}^T \tilde{y}$$

be the transformation of \tilde{y} obtained by forming \tilde{R} from \tilde{X} via sequential orthogonal transformations. To obtain LAR/LASSO directions, we need the solution z_k to the normal equations (2.5) for the current active set:

$$\tilde{X}_k^T \tilde{X}_k z_k = \tilde{X}_k^T \check{y}. \quad (3.8)$$

If \tilde{R}_k is the Cholesky factorization of \tilde{X}_k , and \check{y}_k is the corresponding transformation of \check{y} , then (3.8) is equivalent to

$$\tilde{R}_k z_k = \check{y}_k. \quad (3.9)$$

Use of (3.9) has numerical advantages over (3.8), since growth in roundoff error is bounded by the square root of the condition number of $\tilde{X}_k^T \tilde{X}_k$ rather than the condition number.² The scaled correlations (2.6) are given by

$$c_k = \mathcal{D}_X^{-1} \tilde{X}^T r_k = \mathcal{D}_X^{-1} \tilde{X}^T (\tilde{y} - \tilde{X}b_k) = \mathcal{D}_X^{-1} \tilde{R}^T (\check{y} - \tilde{R}[b_{k-1} + \lambda_k d_k]). \quad (3.10)$$

Step lengths for LAR and LASSO are computed in reduced dimensions from (3.10) instead of (2.7).

Because \tilde{X} and \tilde{R} are related through orthogonal transformation of columns, and there is a one-to-one correspondence between the columns of \tilde{R} and the columns of \tilde{X} ,

$$\|\tilde{R}_{*j}\|_2 = \|\tilde{X}_{*j}\|_2, \quad j = 1, \dots, p,$$

where \tilde{R}_{*j} and \tilde{X}_{*j} , represent the j th column of \tilde{R} and \tilde{X} , respectively, so that

$$\mathcal{D}_X = \text{diag} \left(\|\tilde{R}_{*1}\|_2, \|\tilde{R}_{*2}\|_2, \dots, \|\tilde{R}_{*p}\|_2 \right).$$

Our methods compute \tilde{R} in one pass through the data, which also gives us an efficient way to compute \mathcal{D}_X (which is needed for scaling the correlations) at the outset.³

²The condition number of a positive semi-definite matrix is its largest eigenvalue divided by its smallest eigenvalue. This ratio goes to infinity as the matrix nears singularity.

³Alternatively, R can be scaled initially by its column norms and the coefficients correspondingly transformed after the procedure is completed.

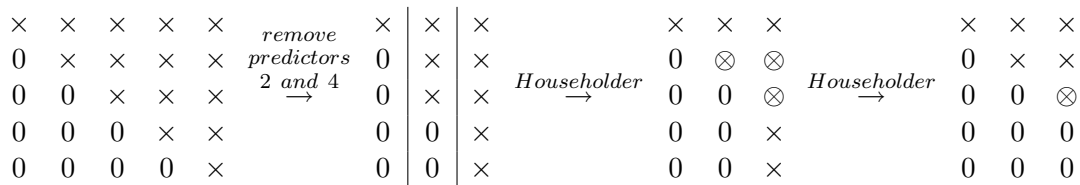


Figure 1: Restoration of triangular form via Householder transformations after removal of some predictors. predictors are kept in their original order. Each arrow corresponds to application of a Householder transformation, \times represents a potentially nonzero entry, and \otimes represents an entry that would typically change value after the Householder transformation is applied.

3.3 Householder Transformations to Introduce Sparsity

A Householder transformation \mathcal{H} is an elementary orthogonal matrix of the form

$$\mathcal{H}(h) \equiv \left(I - \frac{2hh^T}{h^T h} \right),$$

where I is a $d \times d$ matrix, and h is any d vector. A Householder transformation can be applied to change the sparsity structure of a vector or matrix. In particular, if

$$h \equiv \left(\left\| \omega \pm \begin{pmatrix} \omega \\ w \end{pmatrix} \right\|_2 \right), \quad (3.11)$$

where ω is a scalar, then

$$\mathcal{H}(h) \begin{pmatrix} \omega \\ w \end{pmatrix} = \begin{pmatrix} \tilde{\omega} \\ 0 \end{pmatrix} \quad (3.12)$$

(e.g. Golub and Van Loan, section 5.1.3). Because the Householder transformation is orthogonal, the Euclidean norm of vectors to which it is applied remains unchanged. Moreover, vector elements or matrix rows corresponding to zeros in the Householder vector h are unaffected by application of a Householder transformation, so that the transformation can be applied to selected rows of a matrix without affecting the other rows. Details of the formation and application Householder transformations are covered extensively in the literature, see e.g Golub and Van Loan (1996), Chapter 5 and the references cited there.

3.4 Cholesky Factor and Transformed Response for Reduced Sets of Predictors

The Cholesky factor \tilde{R}_k of a reduced set of predictors \tilde{X}_k can be obtained from the Cholesky factor \tilde{R} of the full set of predictors \tilde{X} via orthogonal Householder transformations as illustrated in Figure 1. The dimension of the transformation need only be as large as the number of affected rows, and is essentially the repeated application of relations (3.11) and (3.12) for formation and application of Householder transformations to the parts of the matrix that need to be restored to triangular form. These Householder transformations must also be applied to the vector \tilde{y} to obtain the vector \tilde{y}_k needed in to compute the direction. Starting with \tilde{R} at each stage allows reduction to the Cholesky factor for any subset of predictors, so that predictors dropped at an earlier stage can be reintroduced at later stages, which may be required in the LASSO method.

Although it is not necessary, we assume that the columns of \tilde{R}_k appear in the same order as they would in \tilde{R} , and that the reduced Cholesky factor is formed from the original accumulated Cholesky factor, rather than the Cholesky factor used in the $(k-1)$ st step. However, if at stage k , the active set consists only of leading columns of \tilde{R} or of \tilde{R}_{k-1} , then no update is necessary since \tilde{R}_k and \tilde{y}_k would consist of the corresponding leading columns of \tilde{R} or \tilde{R}_{k-1} and the corresponding leading elements of \tilde{y} or \tilde{y}_{k-1} , respectively. For further discussion on modifying QR factorizations to add or remove columns, see e.g. Chapter 24 of Lawson and Hanson (1995) and Chapter 3 of Miller (2002).

3.5 Detecting Linear Dependence

The approach to linear models and LAR/LASSO via orthogonal factorization has an advantage over the normal equations in terms of detection of ill-conditioning and linear dependence in predictors, which is typically revealed by small diagonal elements in the resulting upper-triangular factor.⁴ When redundant predictors are detected, they can be eliminated using the updating process described in Section 3.4. In the case of LASSO, there should be the option to attempt to reintroduce predictors that are dropped due to ill-conditioning if one or more of the leading predictors present when the ill-conditioning was detected has been dropped from the active set.

3.6 Residual Sum of Squares

Values of the residual sum of squares $\left\| y - \begin{pmatrix} \mathbf{1} & X \end{pmatrix} \begin{pmatrix} \beta_k \\ b_k \end{pmatrix} \right\|_2^2$ are needed for regression diagnostics, such as the C_p statistic proposed in Efron et al. (2004) for variable selection. The residual sum of squares is the square of the Euclidean norm of the projection of y onto the null space of the columns of $\begin{pmatrix} \mathbf{1} & X \end{pmatrix}$ that correspond to nonzero components of b_k .

The residual sum of squares RSS^* for the case in which all components of b_k are nonzero can be obtained during the initial factorization stage: it is the sum of the residual sums of squares for the sequence of regressions defined as the blocks are accumulated. For a given stage k , the corresponding residual sum of squares is then given by

$$RSS_k = RSS^* + \left\| Q^T y - R \begin{pmatrix} \beta_k \\ b_k \end{pmatrix} \right\|_2^2.$$

3.7 Obtaining R and \tilde{R}

In the appendix, we show that if R is a Cholesky factor of $\begin{pmatrix} \mathbf{1} & X \end{pmatrix}$, then a Cholesky factor of \tilde{X} (the matrix \tilde{R} of Section 3) is available as a submatrix of R . We also show that the transformed response $\tilde{y} = \tilde{Q}^T y$ of Section 3 is a subvector of $Q^T y$. Hence orthogonal factorization of $\begin{pmatrix} \mathbf{1} & X \end{pmatrix}$ and simultaneous transformation of y suffices both for computing LAR/LASSO coefficients and for computation of the associated residual sums of squares.

⁴For example, $|r_{jj}| \leq \sqrt{\epsilon_M} \max\{1, \max_k |r_{kk}|\}$ works well in practice as an indicator of linear dependence for the j th predictor, where ϵ_M is the relative machine precision (ϵ_M has the value $2.220446\text{e-}16$ on IEEE compliant machines).

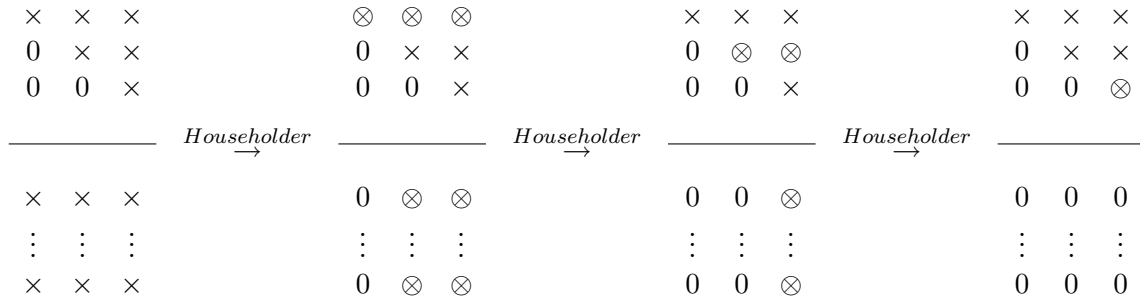


Figure 2: An illustration of the changes in sparsity structure in applying successive Householder transformations in updating the current upper triangular factor (above line) from a new block of rows of predictors (below line). The triangular factor (above the horizontal line) and block (below the horizontal line) need not reside in the same array. Each arrow corresponds to application of a Householder transformation, \times represents a potentially nonzero entry, and \otimes represents an entry that would typically change value after the Householder transformation is applied. In general, k Householder transformations are needed to process a block with k columns.

4 Blockwise Cholesky Factorization

The first step in LAR/LASSO for large data sets is to form the upper triangular Cholesky factor R of $\begin{pmatrix} \mathbf{1} & X \end{pmatrix}^T \begin{pmatrix} \mathbf{1} & X \end{pmatrix}$ in one row-wise pass through the data X . Only a limited number of rows of X will be available at a time, so R is accumulated by applying a succession of orthogonal Householder transformations using the factorization accumulated so far and the currently available rows of X . The basic procedure is known and has been used for large-scale linear regression (e.g. Eldén 2007, Chapter 5), so we give only a brief description here.

As mentioned in Section 3.7, we show in the appendix that it suffices to form the Cholesky factorization of $\begin{pmatrix} \mathbf{1} & X \end{pmatrix}$ via orthogonal transformation and the corresponding transformation of y to obtain the triangular factor \tilde{R} and the transformed response \tilde{y} of Section 3. The corresponding entries of the response y must be available as the rows of X are processed in order to accumulate \tilde{y} , the transformation of y needed for the LAR and LASSO procedures described in Section 3. We assume the data $X \mid y$ is brought into memory in sequential blocks of rows. The procedure as applied to a block of data is illustrated in Figure 2.

If $\begin{pmatrix} \mathbf{1} & X \end{pmatrix}$ resided in memory, its Cholesky factor would be formed in place in the leading $(p+1) \times (p+1)$ submatrix after p Householder transformations, one for each of the leading p columns of the matrix (e.g. Lawson and Hanson 1995, Chapter 3). In our case, we do not have column-wise access to the matrix, and the accumulated upper-triangular matrix R is held separately in memory from the blocks of rows of X as they are accessed. The rowwise formulation of the QR factorization is well known for in-memory data (e.g. Lawson and Hanson 1995, Chapter 27), and has also been extended to data that is sequentially blockwise accessible (e.g. Eldén 2007, Chapter 5). The separation in memory of the current row of the accumulated factor and the current block read in from memory conforms to the natural partition of the Householder vector into its first element and remaining elements (see Section 3.3), so that the accumulated factor need not be combined in a single matrix with the current block in order to achieve the update.

The effect of an individual Householder transformation is as follows: a column of the new block of data is zeroed out, and the value of the corresponding diagonal element of the accumulated factor is changed. The remaining columns of the block, and the corresponding row of the accumulated

factor are also changed. The diagonal elements would correspond to ω (before transformation) and $\tilde{\omega}$ (after transformation) in (3.11) and (3.12), and the associated columns of the block would correspond to w in those relations.

5 Summary and Discussion

A Google Scholar search in April 2007 shows over 400 citations of Efron et al. (2004), and over 950 citations of Tibshirani (1996). The work referencing these papers deals with issues such as categorical and grouped variables, extension to nonlinear models, and cases in which there are many more predictors than observations, but not with datasets with very large numbers of observations. For a recent review, see Hesterberg et al. 2008.

Our approach is an extension to least-angle regression (LAR) and the LASSO method for large datasets that reside out of memory but for which there is sequential access to blocks of rows, as would be the case with standard commercial databases. Our strategy involves a combination sequential blockwise orthogonal transformation techniques that have been used for large-scale linear regression (e.g. Eldén 2007), with methods for updating the QR factorization to accommodate variable sets of predictors (e.g. Lawson and Hanson 1995; Miller 2002). We also show how the correlations for selecting predictors and other quantities needed for the LAR and LASSO methods can be obtained using these factorizations and transformations without accessing the out-of-memory data.

Although our methods do not apply to adaptive linear and nonlinear models because iterative evaluation of functions of the predictors would be required, our approach allows considerable flexibility, since it applies to any model of the form

$$y = (\mathbf{1} \quad \phi_1(X) \quad \dots \quad \phi_m(X)) \alpha + \epsilon,$$

where ϕ_1, \dots, ϕ_m are (possibly nonlinear) functions of the predictors, m is of order p , and the functions are fixed in advance. Such functions could include higher-order terms and interactions between predictors, for example.

Fan et al. (2007) and Fan and Cheng (2007) also propose a blockwise approach to regression and variable selection for sequential blockwise-accessible massive datasets, but they apply their methods to each block separately. Their approach requires synthesis of the blockwise results, and more than one pass is needed to determine an appropriate block size. In our proposed methodology, the results are independent (except for roundoff error) of the order of the observations as well as of the block size, which is limited only by the available memory.

For information on an available S-PLUS implementation of our methods that uses the `bigdata` library to handle datasets that reside out of memory by blockwise processing, as well as on ongoing work on software for generalized least-angle regression, see

<http://www.insightful.com/lars>

References

- Demmel, J. W. (1997) *Applied Numerical Linear Algebra*. SIAM, Philadelphia.
- Efron, B. and Hastie, T. (2003) *LARS software for R and Splus*.
<http://www-stat.stanford.edu/~hastie/Papers/LARS>.
- Efron, B., Hastie, T., Johnstone, I. and Tibshirani, R. (2004) Least angle regression. *Annals of Statistics*, **32**, 407–451.

- Eldén, L. (2007) *Matrix Methods for Data Mining and Pattern Recognition*. SIAM, Philadelphia.
- Fan, T.-H. and Cheng, K.-F. (2007) Tests and variables selection for regression analysis for massive datasets. *Data Knowledge and Engineering*, **63**, 809–817.
- Fan, T.-H., Lin, D. K. J. and Cheng, K.-F. (2007) Regression analysis for massive datasets. *Data Knowledge and Engineering*, **61**, 554–562.
- Golub, G. H. and Van Loan, C. F. (1996) *Matrix Computations*. Johns Hopkins University Press, 3rd edn.
- Hesterberg, T., Choi, N. H., Meier, L. and Fraley, C. (2008) Least angle and l_1 penalized regression: A review. Tech. rep., Insightful Corporation.
- Lawson, C. L. and Hanson, R. J. (1995) *Solving Least Squares Problems*. SIAM, Philadelphia, 2nd edn.
- Leader, J. J. (2004) *Numerical Analysis and Scientific Computation*. Addison Wesley.
- Meyer, C. D. (2000) *Matrix Analysis and Applied Linear Algebra*. SIAM, Philadelphia.
- Miller, A. (2002) *Subset Selection in Regression*. Chapman & Hall, 2nd edn.
- Osborne, M. R., Presnell, B. and Turlach, B. A. (2000) On the LASSO and its dual. *Journal of Computational and Graphical Statistics*, **9**, 319–337.
- Tibshirani, R. (1996) Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society, Series B*, **58**, 267–288.
- Trefethen, L. N. and Bau, D. (1997) *Numerical Linear Algebra*. SIAM, Philadelphia.

Acknowledgments

This work was supported by NIH under NIH SBIR Phase I 1R43GM074313-01 and Phase II 2R44GM074313-02 awards. We thank Tatiana Maravina for reading the manuscript and correcting a number of errors.

A Derivations

In this section, we derive some of the relations on which we have based our version of least-angle regression and LASSO that applies to sequential blockwise accessible out of memory datasets.

Let X be an $n \times p$ matrix, and y and n -vector. Consider the following regression problem:

$$y = \begin{pmatrix} \mathbf{1} & X \end{pmatrix} \begin{pmatrix} \beta \\ b \end{pmatrix} + \epsilon \quad (\text{A.13})$$

where $\mathbf{1}$ is the vector of length n in which all components are equal to 1.

If $\mathbf{b} = \begin{pmatrix} \beta \\ b \end{pmatrix}$ are ordinary least-squares regression coefficients for (A.13), then b are ordinary least-squares regression coefficients for

$$\tilde{y} + \tilde{X}b + \epsilon, \quad (\text{A.14})$$

where \tilde{y} represents y with its mean subtracted out

$$\tilde{y} \equiv y - \bar{y}\mathbf{1},$$

and \tilde{X} represents X with column means subtracted out, that is,

$$\tilde{X} \equiv X - \frac{\mathbf{1}\mathbf{1}^T}{n}X = \left(I - \frac{\mathbf{1}\mathbf{1}^T}{n} \right) X,$$

and

$$\beta = \frac{\mathbf{1}^T(y - Xb)}{n} = \bar{y} - (\bar{X}_{*1} \quad \bar{X}_{*2} \quad \dots \quad \bar{X}_{*p})^T b, \quad (\text{A.15})$$

where \bar{X}_{*j} is the mean of the j th column of X .

Suppose that we have a QR factorization of $\begin{pmatrix} \mathbf{1} & X \end{pmatrix}$:

$$\begin{pmatrix} \mathbf{1} & X \end{pmatrix} = \mathcal{Q}\mathcal{R} = \mathcal{Q} \begin{pmatrix} R \\ O \end{pmatrix}, \quad (\text{A.16})$$

where \mathcal{Q} is an $n \times n$ matrix with orthogonal columns, and \mathcal{R} is an $n \times (p+1)$ matrix, and R is a $(p+1) \times (p+1)$ upper triangular matrix. Then, if

$$\mathcal{Q} = \begin{pmatrix} Q & Z \end{pmatrix},$$

where Q is $n \times (p+1)$, an alternative QR factorization is

$$\begin{pmatrix} \mathbf{1} & X \end{pmatrix} = QR \quad (\text{A.17})$$

(see Lawson and Hanson, Chapters 2–3 or Golub and Van Loan 1996, Chapter 5).

For a more abstract mathematical treatment of relevant concepts in linear algebra, see Meyer (2000), Chapter 5. Let R have the following partition

$$R = \begin{pmatrix} \rho & s^T \\ 0 & U \end{pmatrix}. \quad (\text{A.18})$$

In **Proposition 1**, we derive the relationship between Cholesky factors of $\begin{pmatrix} \mathbf{1} & X \end{pmatrix}$ and \tilde{X} . In **Proposition 2**, we derive a relationship between the transformations of y corresponding to a QR factorization of $\begin{pmatrix} \mathbf{1} & X \end{pmatrix}$ and the transformation of \tilde{y} corresponding to a QR factorization of \tilde{X} . In **Proposition 3**, we derive an expression for the regression coefficient of the intercept in (A.13) in terms of the QR factorization of $\begin{pmatrix} \mathbf{1} & X \end{pmatrix}$ and the corresponding transformation of y .

Proposition 1: Suppose that R is an upper triangular Cholesky factor $\begin{pmatrix} \mathbf{1} & X \end{pmatrix}^T \begin{pmatrix} \mathbf{1} & X \end{pmatrix}$:

$$R^T R = \begin{pmatrix} \mathbf{1} & X \end{pmatrix}^T \begin{pmatrix} \mathbf{1} & X \end{pmatrix},$$

R $p \times p$ upper triangular. Let R be partitioned as in (A.18). Then the upper triangular matrix U in (A.18) is a Cholesky factor of \tilde{X} .

Proof.

$$R^T R = \begin{pmatrix} \rho & s^T \\ 0 & U \end{pmatrix}^T \begin{pmatrix} \rho & s^T \\ 0 & U \end{pmatrix} = \begin{pmatrix} \rho^2 & \rho s^T \\ \rho s & U^T U + s s^T \end{pmatrix}, \quad (\text{A.19})$$

and

$$\begin{pmatrix} \mathbf{1} & X \end{pmatrix}^T \begin{pmatrix} \mathbf{1} & X \end{pmatrix} = \begin{pmatrix} \mathbf{1}^T \\ X^T \end{pmatrix} \begin{pmatrix} \mathbf{1} & X \end{pmatrix} = \begin{pmatrix} \mathbf{1}^T \mathbf{1} & \mathbf{1}^T X \\ X^T \mathbf{1} & X^T X \end{pmatrix} = \begin{pmatrix} n & \mathbf{1}^T X \\ X^T \mathbf{1} & X^T X \end{pmatrix}.$$

Since $R^T R = \begin{pmatrix} \mathbf{1} & X \end{pmatrix}^T \begin{pmatrix} \mathbf{1} & X \end{pmatrix}$, equating matrix partitions gives

$$\begin{aligned} \rho^2 &= n \\ X^T \mathbf{1} &= \rho s \\ X^T X &= U^T U + s s^T, \end{aligned} \quad (\text{A.20})$$

from which it follows that

$$\begin{aligned} U^T U &= X^T X - s s^T = X^T X - \frac{X^T \mathbf{1} \mathbf{1}^T X}{\rho^2} = X^T X - \frac{X^T \mathbf{1} \mathbf{1}^T X}{n} \\ &= X^T \left(I - \frac{\mathbf{1} \mathbf{1}^T}{n} \right) X = X^T \left(I - \frac{\mathbf{1} \mathbf{1}^T}{n} \right) \left(I - \frac{\mathbf{1} \mathbf{1}^T}{n} \right) X = \tilde{X}^T \tilde{X}, \end{aligned}$$

so that U is a Cholesky factor of \tilde{X} .

Proposition 2: Let y be any n vector, and consider the QR factorization (A.17). Let R be partitioned as in (A.18). and partition $Q^T y$ as

$$Q^T y = \begin{pmatrix} \zeta \\ z \end{pmatrix},$$

where ζ is a scalar. Then z differs from the transformation $\tilde{y} = \tilde{Q}^T \tilde{y}$ only in the null space of U^T . Moreover, if U is nonsingular (or, equivalently, the columns of \tilde{X} are linearly independent), $z = \tilde{y}$.

Proof. Let $\mathbf{b} = \begin{pmatrix} \beta \\ b \end{pmatrix}$ be ordinary least-squares regression coefficients for (A.13). Then,

$$\begin{aligned} \begin{pmatrix} \mathbf{1} & X \end{pmatrix}^T \begin{pmatrix} \mathbf{1} & X \end{pmatrix} \mathbf{b} &= \begin{pmatrix} n & \mathbf{1}^T X \\ X^T \mathbf{1} & X^T X \end{pmatrix} \begin{pmatrix} \beta \\ b \end{pmatrix} = \begin{pmatrix} n\beta + \mathbf{1}^T X b \\ \beta X^T \mathbf{1} + X^T X b \end{pmatrix} = \begin{pmatrix} \mathbf{1} & X \end{pmatrix}^T y \\ &= R^T Q^T y = \begin{pmatrix} \rho & s^T \\ 0 & U \end{pmatrix}^T \begin{pmatrix} \zeta \\ z \end{pmatrix} = \begin{pmatrix} \rho \zeta \\ \zeta w + U^T z \end{pmatrix}. \end{aligned}$$

Using relations from (A.20), this can be rewritten as

$$\begin{pmatrix} n\beta + \mathbf{1}^T X b \\ \beta X^T \mathbf{1} + X^T X b \end{pmatrix} = \begin{pmatrix} n\beta + \mathbf{1}^T X b \\ \beta X^T \mathbf{1} + \left[U^T U + \frac{X^T \mathbf{1} \mathbf{1}^T X}{n} \right] b \end{pmatrix} = \begin{pmatrix} \sqrt{n} \zeta \\ \zeta \frac{X^T \mathbf{1}}{\sqrt{n}} + U^T z \end{pmatrix}. \quad (\text{A.21})$$

So that

$$\zeta = \frac{n\beta + \mathbf{1}^T Xb}{\sqrt{n}},$$

and

$$\zeta \frac{X^T \mathbf{1}}{\sqrt{n}} = \left[\frac{n\beta + \mathbf{1}^T Xb}{\sqrt{n}} \right] \frac{X^T \mathbf{1}}{\sqrt{n}} = \beta X^T \mathbf{1} + \frac{(\mathbf{1}^T Xb)(X^T \mathbf{1})}{n} = \beta X^T \mathbf{1} + \frac{(X^T \mathbf{1})(\mathbf{1}^T Xb)}{n} = \beta X^T \mathbf{1} + \frac{X^T \mathbf{1} \mathbf{1}^T X}{n} b.$$

Combining with (A.21), this gives

$$U^T U b = U^T z.$$

Now U is a Cholesky factor of \tilde{X} from **Proposition 2** and b is an ordinary least-squares solution of (A.14), so

$$U^T U b = \tilde{X}^T \tilde{y} = U^T \tilde{Q}^T \tilde{y} = U^T \check{y},$$

by the discussion in Section 3.1. Hence

$$U^T z = U^T \check{y} \quad \text{or} \quad U^T (z - \check{y}) = 0,$$

so that z and \check{y} differ only in the null space of U^T . Moreover, if U is nonsingular,

$$z = U^{-T} U^T z = U^{-T} U^T \check{y} = \check{y}.$$

Proposition 3: Let $\mathbf{b} = \begin{pmatrix} \beta \\ b \end{pmatrix}$ be ordinary least-squares regression coefficients for (A.13). Let R in the QR factorization (A.16) be partitioned as in (A.18), and let $\mathcal{Q}^T y$ be partitioned as

$$\mathcal{Q}^T y = \begin{pmatrix} Q^T y \\ Z^T y \end{pmatrix} = \begin{pmatrix} \zeta \\ z \\ Z^T y \end{pmatrix},$$

where ζ is a scalar. Then

$$\beta = \frac{\rho(\zeta - s^T b)}{n}.$$

Proof. From (A.16), we have

$$\mathcal{Q}^T \begin{pmatrix} \mathbf{1} & X \end{pmatrix} = \begin{pmatrix} \mathcal{Q}^T \mathbf{1} & \mathcal{Q}^T X \end{pmatrix} \begin{pmatrix} R \\ O \end{pmatrix}.$$

From the partition (A.18), we have

$$\mathcal{Q}^T \mathbf{1} = \begin{pmatrix} \rho \\ 0 \\ 0 \end{pmatrix} \quad \text{and} \quad \mathcal{Q}^T X = \begin{pmatrix} s^T \\ U \\ O \end{pmatrix}.$$

Combining these results with (A.15), we have

$$\beta = \frac{\mathbf{1}^T (y - Xb)}{n} = \frac{\mathbf{1}^T \mathcal{Q} \mathcal{Q}^T (y - Xb)}{n} = \frac{1}{n} \begin{pmatrix} \rho \\ 0 \\ 0 \end{pmatrix}^T \left[\begin{pmatrix} \zeta \\ z \\ Z^T y \end{pmatrix} - \begin{pmatrix} s^T \\ U \\ O \end{pmatrix} b \right] = \frac{\rho(\zeta - s^T b)}{n}.$$