

# Sign-constrained least squares estimation for high-dimensional regression

Nicolai Meinshausen  
University of Oxford, UK  
meinshausen@stats.ox.ac.uk

January 24, 2012

## Abstract

Many regularization schemes for high-dimensional regression have been put forward. Most require the choice of a tuning parameter, using model selection criteria or cross-validation schemes. We show that a simple non-negative or sign-constrained least squares is a very simple and effective regularization technique for a certain class of high-dimensional regression problems. The sign constraint has to be derived via prior knowledge or an initial estimator but no further tuning or cross-validation is necessary. The success depends on conditions that are easy to check in practice. A sufficient condition for our results is that most variables with the same sign constraint are positively correlated. For a sparse optimal predictor, a non-asymptotic bound on the L1-error of the regression coefficients is then proven. Without using any further regularization, the regression vector can be estimated consistently as long as  $\log(p)s/n \rightarrow 0$  for  $n \rightarrow \infty$ , where  $s$  is the sparsity of the optimal regression vector,  $p$  the number of variables and  $n$  sample size. Network tomography is shown to be an application where the necessary conditions for success of non-negative least squares are naturally fulfilled and empirical results confirm the effectiveness of the sign constraint for sparse recovery.

## 1 Introduction

High-dimensional regression problems are characterized by a large number of predictor variables in relation to sample size. Regularization (in a broad sense) is of critical importance for high-dimensional problems and much attention has been paid to various schemes and their properties in recent years, including the *Ridge* estimator [Hoerl and Kennard, 1970], *non-negative Garrote* [Breiman, 1995], the *Lasso* [Tibshirani, 1996] and various variations of the latter, including the *group Lasso* [Yuan and Lin, 2006] and *adaptive Lasso* [Zou, 2006]. Datasets with very low signal-to-noise ratio offer similar challenges to high-dimensional problems even if the notional sample size is quite high.

Sign-constraints on the regression coefficients are a simpler regularization and have been first advocated by I.J. Good, as covered in the book Lawson and Hanson [1995]. There is a wide range of problems where the sign of the regression coefficients can either be estimated by an initial estimator or where it is known a priori, such as in image processing and spectral analysis [Waterman, 1977, Bellavia et al., 2006, Donoho et al., 1992, Chen and Plemmons, 2009]. Sign-constraints have also been implemented for matrix factorizations, specifically the *non-negative Matrix factorization* [Lee et al., 1999, Lee and Seung, 2001, Ding et al., 2010] and *non-negative least squares* regression can be a useful tool for this factorization [Kim and Park, 2007]. We study the performance of non-negative least squares type problems under a so-called *Positive Eigenvalue Condition*, which can be checked for any given dataset by solving a quadratic programming problem. A sufficient condition uses only the minimum of all entries in the design matrix. It is shown that non-negative (or, in general, sign-constrained) least squares is a surprisingly effective regularization technique for high-dimensional regression problems under these conditions. If the *Positive Eigenvalue Condition* is not fulfilled, the sign constraint is still a good ingredient in a regularization framework. The *non-negative Garrote*

[Breiman, 1995] is, for example, making use of a sign-constraint, where the signs are derived from an initial estimator as is the *positive Lasso* [Efron et al., 2004].

The data are assumed to be given by a  $n \times 1$ -vector of real-valued observations  $\mathbf{Y}$  and a  $n \times p$ -dimensional matrix  $\mathbf{X}$ , where column  $k$  of  $\mathbf{X}$  contains all  $n$  samples of the  $k$ -th predictor variable for  $k = 1, \dots, p$ . The non-negative least squares (NNLS) regression estimator is defined as

$$\hat{\beta} := \operatorname{argmin}_{\beta} \|\mathbf{Y} - \mathbf{X}\beta\|_2^2 \quad \text{such that} \quad \min_k \beta_k \geq 0. \quad (1)$$

We will work with a positivity constraint without limitation of generality since variables that are constrained to be negative can be replaced by their negative counterpart and the problem can thus always be framed as a non-negative least squares optimisation. Problem (1) is a convex optimization problem and can be solved with general quadratic programming problem solvers, including active set [Lawson and Hanson, 1995], iterative [Kim et al., 2006] and interior-point approaches [Bellavia et al., 2006]. A tailor-made fast approximate algorithm based on random projections has recently been proposed in Boutsidis and Drineas [2009]. The recent manuscript Slawski et al. [2011] contains independent work on the behaviour of NNLS in high-dimensions. Using the same *Positive Eigenvalue Condition* (which is called self-regularizing design condition), a bound on the prediction error of NNLS and a sparse recovery property after hard thresholding are shown in Slawski et al. [2011]. Our main focus is on sparse recovery in the  $\ell_1$ -sense. The bounds on prediction error are also of different nature since the assumptions are different. We make use of the so-called compatibility condition which appears in most sparse recovery results in the  $\ell_1$ -norm penalized estimation literature [Van De Geer and Bühlmann, 2009] and derive, with the help of this condition, tight non-asymptotic bounds on the prediction error.

Note that the non-negative least squares estimator (1) does not require the choice of a tuning parameter beyond choosing the sign of the coefficients. Imposing a sign-constraint might seem like a very weak regularization but it will be shown that the estimator is remarkably different from the un-regularized least squares estimator. It can cope with high-dimensional problems, where the number of predictor variables vastly exceeds sample size. It will be shown to be a consistent estimator as long as the underlying optimal prediction is sufficiently sparse (ie using only a small subset of all predictor variables) and the so-called *Positive Eigenvalue Condition* is fulfilled.

The manuscript is organized as follows. The notation and the main two assumptions, the compatibility and *Positive Eigenvalue Condition*, are introduced in Section 2. Our main result, a  $\ell_1$ -bound on the difference between the NNLS estimator and the optimal regression coefficients, is shown in Section 3, along with a bound on the prediction error.

## 2 Notation and Assumptions

We assume that the  $n$  samples  $\mathbf{Y} \in \mathbb{R}^n$  are drawn from  $\mathbf{X}\beta^* + \varepsilon$  for some  $p$ -dimensional vector  $\beta^*$  with  $\min_k \beta_k^* \geq 0$  and  $\varepsilon \sim \mathcal{N}(0, \sigma^2)$  for some  $\sigma > 0$ . Let  $S$  be the set of non-zero entries of the optimal solutions,  $S = \{k : \beta_k^* \neq 0\}$  and  $N = S^c$  be the complement of  $S$ . We could also let  $\beta^*$  be the best approximation to the data-generating model under positivity constraints but will refrain from doing so for notational simplicity. We assume that the columns of  $\mathbf{X}$  are standardized to  $\ell_2$ -norm of  $n$ . Despite not necessarily assuming that the columns are mean-centered, we call  $\hat{\Sigma} = n^{-1}\mathbf{X}^T\mathbf{X}$  the covariance matrix throughout.

We make two major assumptions for the main result, one about sparse eigenvalues and another about the positive eigenvalue between predictor variables.

### 2.1 Compatibility Condition

There has been much recent work on the properties of the Lasso [Tibshirani, 1996]. Many similar conditions for success of the Lasso penalization schemes have been derived [for example Zhang and Huang, 2008, Meinshausen and Yu, 2009, Wainwright, 2009, Bunea et al., 2007, 2006, Van De Geer, 2008, Bickel et al., 2009]. A good overview of all conditions and their relations is given in Van De Geer and Bühlmann [2009]. The weakest condition is based on the notion of  $(L, S)$  restricted  $\ell_1$ -eigenvalues.

The  $(L, S)$  restricted  $\ell_1$ -eigenvalue of matrix  $\mathbf{A}$  is defined as:

$$\phi_{\text{compatible}}^2(L, S, \mathbf{A}) := \min \left\{ s \frac{\beta^T \mathbf{A} \beta}{\|\beta\|_1^2} : \beta \in \mathcal{R}(L, S) \right\},$$

where  $\mathcal{R}(L, S) = \{\beta : \|\beta_N\|_1 \leq L\|\beta_S\|_1\}$  and  $s = |S|$ .

A lower bound on this restricted eigenvalue is necessary for success of the Lasso, either in a prediction loss or coefficient recovery sense and was called the compatibility condition in Van De Geer and Bühlmann [2009]. It was shown to be weaker than all similar conditions such as the *Restricted Isometry Property* [Candes and Tao, 2007].

We make the following assumption.

**Assumption 1** (Compatibility Condition). *There exists some  $\phi > 0$  such that the  $(L, S)$ -restricted  $\ell_1$ -eigenvalue  $\phi_{\text{compatible}}^2(L, S, \hat{\Sigma}) \geq \phi$ .*

The value of  $L$  will be specified in Theorem 1.

**Remark 1.** *The assumption is formulated for the empirical covariance matrix  $\hat{\Sigma}$  but can also easily be reformulated on the population covariance matrix  $\Sigma$  for random design. Assume that the maximal difference between the population and empirical covariance matrix is bounded by  $\delta > 0$ , that is  $\|\hat{\Sigma} - \Sigma\|_\infty \leq \delta$ . This assumption is fulfilled with high probability for many data sets with larger sample size. If the predictors have for example a multivariate normal distribution (which will not be assumed elsewhere), then the condition is fulfilled with probability  $1 - 2\exp(-t)$  for  $\delta \geq \sqrt{u} + u$  with  $u = (4t + 8\log(p))/n$ , see (10.1) in Van De Geer and Bühlmann [2009]. If  $\delta \leq \phi^2 / (4(L + 1)^2 s)$ , then  $\phi_{\text{compatible}}^2(L, S, \Sigma) \geq \phi$  implies  $\phi_{\text{compatible}}^2(L, S, \hat{\Sigma}) \geq \phi/2$ . The proof follows from the inequality  $\phi_{\text{compatible}}^2(L, S, \hat{\Sigma}) \geq \phi_{\text{compatible}}^2(L, S, \Sigma) - (L + 1)\sqrt{\delta}s$  in Corollary 10.1 in Van De Geer and Bühlmann [2009]. The Compatibility Condition could thus be imposed on the population covariance matrix instead of the empirical covariance matrix.*

## 2.2 Positive eigenvalue condition

The following *Positive Correlation Condition* is the main assumption necessary to show success of non-negative least squares.

The positively constrained minimal  $\ell_1$ - eigenvalue of matrix  $\mathbf{A}$  is defined as

$$\phi_{\text{pos}}^2(\mathbf{A}) := \min \left\{ \frac{\beta^T \mathbf{A} \beta}{\|\beta\|_1^2} : \min_k \beta_k \geq 0 \right\},$$

A lower bound on this restricted eigenvalue will be a sufficient condition for sparse recovery success of NNLS.

**Assumption 2** (Positive Eigenvalue Condition). *There exists some  $\nu > 0$  such that  $\phi_{\text{pos}}^2(\hat{\Sigma}) \geq \nu$ .*

A lower bound on this eigenvalue seems to be a much stricter condition than the *Compatibility Condition*. However, the latter allows for positive and negative regression coefficients, while the *Positive Eigenvalue Condition* is restricted to positive coefficients. There are thus some immediate examples where it is fulfilled, which we discuss below.

**Example I: strictly positive covariance matrix.** The *Positive Correlation Condition* is fulfilled if  $\min_{i,j} \hat{\Sigma}_{ij} \geq \nu > 0$ , that is all entries in the covariance matrix are strictly positive. Again, this condition could also be formulated for the population covariance matrix, using a bound on  $\|\Sigma - \hat{\Sigma}\|_\infty$ .

We also remark on the case of general sign-constraints (some variables constrained to be positive, some negative). The condition applies then to the dataset where all variables with a negativity constraint have been replaced with their negative counterparts. The constraint on the original covariance matrix is thus that it forms two blocks. The variables in the first block are the variables with a positivity constraint and the second block is formed by all variables with a negativity constraint. Correlations are required to be positive within a block and negative between blocks.

A generalization of Example I is the following.

**Example II: only few negative entries.** Let  $\mathcal{A} := \{i : \hat{\Sigma}_{ij} < 0 \text{ for some } 1 \leq j \leq p\}$  be the minimal set such that  $\hat{\Sigma}_{ij} < 0$  implies  $\{i, j\} \subseteq \mathcal{A}$  for all  $1 \leq i, j \leq p$ . The *Positive Eigenvalue Condition* is fulfilled if both of the conditions below are fulfilled for some  $\nu > 0$ .

1. All entries of the covariance matrix are strictly positive on  $\mathcal{A}^c$ , that is  $\hat{\Sigma}_{ij} \geq 2\nu$  if  $\{i, j\} \subseteq \mathcal{A}^c$  for all  $1 \leq i, j \leq n$ .
2. A restricted eigenvalue condition holds on the set  $\mathcal{A}$ , ie

$$\min \left\{ \frac{\beta^T \hat{\Sigma} \beta}{\|\beta\|_1^2} : \beta_k = 0 \text{ for all } k \in \mathcal{A}^c \right\} > 2\nu.$$

If the set  $\mathcal{A}$  is very small, in particular much smaller than  $n$ , the latter restricted  $\ell_1$ -eigenvalue condition is in general not very restrictive. The important criterion is thus whether the set  $\mathcal{A}$  is small compared to the sample size.

**Example III: block matrix.** For a  $p \times p$ -matrix  $\mathbf{A}$  and a set  $K \subseteq \{1, \dots, p\}$ , let  $\mathbf{A}_{KK}$  be the  $|K| \times |K|$ -submatrix formed by all elements in set  $K$ . Suppose

1. Entries of the covariance matrix can be negative but fulfil  $\hat{\Sigma}_{ij} \geq -\rho/p^2$  for all  $1 \leq i, j \leq n$  and some  $\rho > 0$ .
2. The set of variables  $\{1, \dots, p\}$  can be partitioned into  $B \geq 1$  blocks  $B_j \subseteq \{1, \dots, p\}$  such that  $\phi_{pos}^2(\hat{\Sigma}_{B_j B_j}) \geq (\nu + \rho)B$  for all  $j = 1, \dots, B$ .

A more specific example is thus: all entries in  $\hat{\Sigma}$  are larger than  $-\rho/p^2$  for some  $\rho > 0$  and  $\hat{\Sigma}_{ij} \geq (\nu + \rho)B$  if both  $i, j$  are within the same block.

The *Positive Eigenvalue Condition* is fulfilled with parameter  $\nu > 0$ .

The positive aspect of the condition is that it is very easy to check in practice whether it applies (at least approximately) and whether one would thus expect the bounds shown below to apply to a given dataset.

### 3 Main Results

It will be shown that non-negative least squares leads to a good recovery of the optimal sparse regression vector for high-dimensional data. We study the  $\ell_1$ -error in the regression vector, which also yields a bound on the  $\ell_2$ -error and prediction loss.

**Theorem 1.** *Assume that the Positive Eigenvalue Condition holds with  $\nu > 0$ . Choose any  $0 < \eta < 1/3$ . Assume that the compatibility condition holds with  $\phi > 0$  for  $L = 4\nu^{-1}$ . Setting*

$$K_{p,\eta}^2 := 2 \log \left( \frac{\sqrt{2p}}{\sqrt{\pi\eta}} \right)$$

and assuming  $\min_{k \in S} \beta_k > K_{p,\eta} \sigma / \sqrt{n\phi}$ , it then holds with probability at least  $1 - \eta$  that

$$\|\hat{\beta} - \beta^*\|_1 \leq K_{p,\eta} (5/\nu + 4/\sqrt{\phi}) \frac{s\sigma}{\sqrt{n}} \quad (2)$$

A proof is in the appendix.

The result might be surprising since it implies that non-negative least squares is succeeding in recovering the regression coefficients in an  $\ell_1$ -sense if  $\log(p)s/\sqrt{n} \rightarrow 0$  for  $n \rightarrow \infty$ , a scaling that requires for general design a lot more regularization in the form of Lasso penalties (or similar).

The result does not imply exact sign recovery in the sense that the non-zero coefficients equal exactly the set  $S$  (and indeed this will in general not be the case), but it implies that the  $S$  largest coefficients correspond to the variables in the set  $S$ .

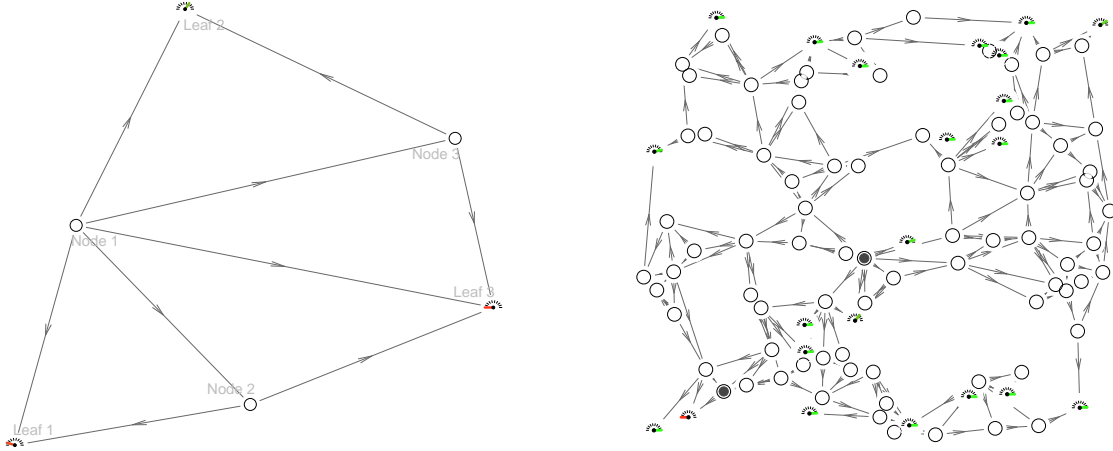


Figure 1: Left: A network with three internal nodes and three leaf nodes. The (unobservable) losses at the internal nodes are  $(10,10,0)$ , meaning that the first two nodes lead to a loss rate of 10 and the third node is not leading to any losses. The observations of the loss rates at the leaf nodes are then  $(8,3,9)$ . Using the observations at the leaf nodes and knowledge of the topology, NNLS can correctly identify the two first nodes as responsible for the losses. Right: A network with 78 internal nodes and 22 leaf nodes. Two of the internal nodes have a positive loss (marked with a dot) and the observations at the leaf nodes are again sufficient to pinpoint the (unknown) location of the two nodes using NNLS estimation.

**Corollary 1.** *Under the same conditions as Theorem 1 and the stronger assumption that the minimum over all non-zero coefficients is bounded from below by  $\min_{k \in S} \beta_k \geq 2K_{p,\eta}\sigma(5/\nu + 4/\sqrt{\phi})s/\sqrt{n}$ , it holds with probability at least  $1 - \eta$  that the indices of the  $s$  largest absolute coefficients in  $\hat{\beta}$  are identical to the set  $S$ .*

This follows immediately from Theorem 1 since the  $\ell_1$ -bound on the difference between  $\hat{\beta}$  and  $\beta^*$  implies the same bound in the supremum-norm.

The bound in Theorem 1 also implies a bound on the prediction error.

**Theorem 2.** *Under the same conditions as Theorem 1, with probability at least  $1 - \eta$  for any  $0 < \eta < 1/3$ ,*

$$\|\mathbf{X}(\hat{\beta}^{oracle} - \hat{\beta})\|_2^2 \leq 2K_{p,\eta}^2\sigma^2(5/\nu + 2/\sqrt{\phi})s.$$

A proof is given in the appendix. The mean squared error, introduced by using NNLS instead of the oracle estimator is thus proportional to  $\log(p)^2s/n$ . The result implies asymptotically vanishing prediction error if  $s \log(p)^2/n \rightarrow 0$  for  $n \rightarrow \infty$ .

## 4 Numerical Results

The results above imply that NNLS can be very effective if (a) the sign of regression coefficients is known or can easily be estimated and (b) the *Positive Eigenvalue Condition* holds. *Network tomography* is a good example [Castro et al., 2004]. For others, including image analysis and applications in signal processing, see Slawski et al. [2011]. There are different aspects of network tomography, including origin-destination matrix

estimation and link-level network tomography; see Castro et al. [2004] for a good overview of the statistical aspects and Xi et al. [2006] and Lawrence et al. [2006] for a discussion of active tomography in the context of link-level analysis. We will focus on one aspect of the link-level network tomography. The network consists of nodes arranged in a directed acyclic graph (or sometimes as a special case a tree) and measurements can be taken at the leaf nodes. These measurements are used to infer the state of all the nodes in the network. In a communication network, the measurements can be the delay or loss rate of packages, in a transport network (such as water distributions networks) it can be the shortfall of the flow rate compared to the expected rate. Since the network topology is assumed to be known, the measurements consist typically of noise plus a linear combination of the internal and unobservable states of the nodes in the network. If a node in the network has a loss (be it in the form of delaying packages or loss of water flow), it will have a linear effect on all leaf nodes that are descendants of the node in the directed acyclic graph.

Figure 1 shows a toy example. Imagining a flow passing through the tree from the internal nodes to the leaf nodes, the entry  $\mathbf{X}_{i,j}$  is the proportion of flow in node  $j$  that reaches leaf node  $i$  if flow is divided equally among all outgoing edges in each node of the tree. Three internal nodes have loss rates  $(\beta_1, \beta_2, \beta_3) = (10, 10, 0)$ . The loss rates  $\mathbf{Y} = (Y_1, Y_2, Y_3)$  at the three leaf nodes are then given by  $\mathbf{Y} = \mathbf{X}\beta + \varepsilon$  for some i.i.d. noise  $\varepsilon$  and

$$\mathbf{X} = \begin{pmatrix} 0.3 & 0.5 & 0 \\ 0.3 & 0 & 0.5 \\ 0.4 & 0.5 & 0.5 \end{pmatrix}.$$

A positivity constraint on the coefficient vectors is clearly appropriate since there will in general not be a negative loss at internal nodes (for example no unexpected *gain* of water in a distribution network). In the noiseless case, the NNLS solution recovers exactly the internal states  $(10, 10, 0)$  and thus identifies correctly the first two nodes as responsible for the loss of the flow rate in all three leaf nodes. In this simple example, the number of leaf nodes is equal to the number of internal nodes and ordinary least squares would also work in the noiseless case. Least squares clearly ceases to be useful once the number of internal nodes exceeds the number of leaf nodes. Note that, contrary to the previous literature (for example Castro et al. [2004], Lawrence et al. [2006]) we do not attempt to fit a stochastic model to the observations. We are merely trying to directly estimate the current internal state  $\beta$  of the nodes in the network as accurately as possible.

The theory suggests that a non-negativity constraint can already be very powerful under certain constraints on the design matrix. The main condition is the *Positive Eigenvalue Condition*. In our simple network tomography example, it is obvious that all entries in  $\mathbf{X}$  are positive and the same is hence true for  $\hat{\Sigma} = n^{-1}\mathbf{X}^T\mathbf{X}$ . Entries in  $\mathbf{X}$  correspond to the amount of loss (delay of packages or reduction in flow rate) in a leaf node caused by a specific loss at an internal node and is non-zero if and only if there is a connection between the internal and the leaf node. Suppose that all non-zero entries in  $\mathbf{X}$  have entries at least as large as  $\delta$  for some  $\delta > 0$ . Suppose further that we can group all internal nodes into  $B$  blocks such that the internal nodes within a block share at least one leaf node to which they all connect. The *Positive Eigenvalue Condition* is then fulfilled with value  $\delta^2/B$ ; see Example III in the discussion of the condition.

The theory seems to show that under these conditions the NNLS-regularization is effective. To test this, we examine the effect of placing an additional  $\ell_1$ -constraint on the coefficient by computing

$$\beta^\lambda := \operatorname{argmin}_\beta \|\mathbf{Y} - \mathbf{X}\beta\|_2^2 \quad \text{such that } \min_k \beta_k \geq 0 \text{ and } \|\beta\|_1 \leq \lambda. \quad (3)$$

Let  $\hat{\beta}$  be again the NNLS-solution defined in (1). It is obvious that  $\beta^\lambda \equiv \hat{\beta}$  for all  $\lambda \geq \lambda_{\max}$  for  $\lambda_{\max} := \|\hat{\beta}\|_1$ .

We generate networks of similar type as the ones shown in Figure 1. The number  $N$  of total nodes is chosen for each of 1000 simulations uniformly out the set  $\{25, 50, 100, 200, 400\}$ . Nodes are distributed uniformly on the area  $[-1, 1]^2$  and numbered in order of their Euclidean distance from the origin. Starting with the first node  $k = 1$  closest to the origin, edges are drawn between it and its  $K$  nearest neighbours with a larger ordering number (where  $K$  is drawn uniformly from the set  $\{5, 10, 20\}$ ). When drawing edges at node  $k = 1, \dots, N - 1$ , they are deleted with probability  $\nu$  (where  $\nu$  is drawn uniformly from the set  $\{.2, .4, .6, .8, 1\}$ ) or when the edge would cross a previously drawn edge. Imagining again a flow passing through the tree from the internal nodes to the leaf nodes, the entry  $\mathbf{X}_{i,j}$  is the proportion of flow in node

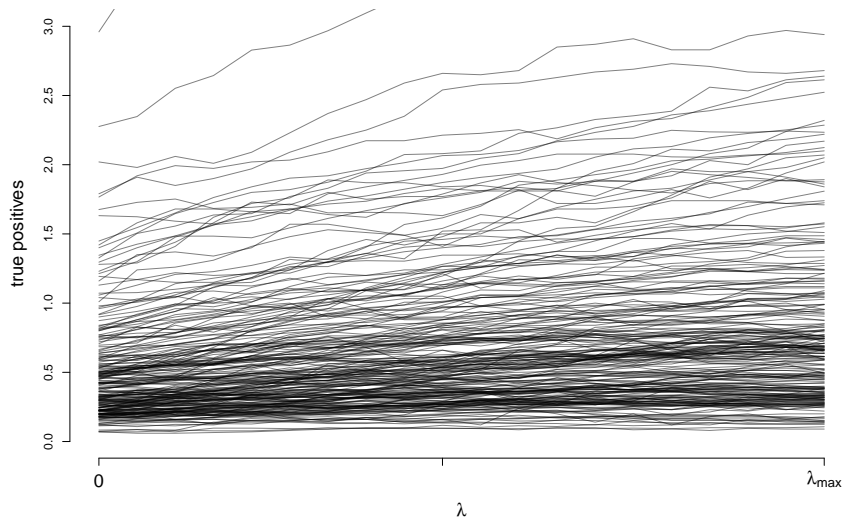


Figure 2: The average number of correctly identified internal nodes with a positive loss under 1000 different scenarios with an additional  $\ell_1$ -constraint as in (3). The NNLS solution corresponds to  $\lambda = \lambda_{\max}$  and is seen to be in general superior to the solutions under additional shrinkage.

$j$  that reaches leaf node  $i$  if flow is divided equally among all outgoing edges in each node of the tree. For each of the 1000 simulations, we draw a single graph from the parameters as described above and also draw the noise variance uniformly from the set  $\{0, 0.125, 0.25, 0.5, 1, 2, 4\}$  and a number  $s$  of non-zero entries in  $\beta$  (corresponding to nodes with a delay or loss), where  $s$  is drawn uniformly from the set  $\{2, 5, 10\}$ . The  $s$  non-zero entries from  $\beta$  are generated independently as the absolute value of a standard-normal random variable. For each such setting, we simulate 50 times the vector  $\mathbf{Y}$  and reconstruct with  $\hat{\beta}^\lambda$  as in (3) for an evenly spaced grid of 20 points between  $\lambda = 0$  and  $\lambda = \lambda_{\max}$ , the NNLS solution. Nodes are put in decreasing order of the reconstructed value  $\hat{\beta}^\lambda$ . We record the first entry in the re-ordered vector  $\hat{\beta}^\lambda$  that corresponds to a false positive (a zero entry in the equally re-ordered vector  $\beta$ ) and call the number of true positives the number of values of  $\hat{\beta}^\lambda$  with larger value than the first false positive.

Figure 2 shows the average number of true positives as a function of  $\lambda$ . Each line corresponds to the average value over all 50 simulations in a given scenario. For nearly all scenarios there is no benefit in placing an additional  $\ell_1$ -penalty on the coefficients. The NNLS solution is thus a very good and simple estimator in these settings, as expected from theory. Additional regularization by an  $\ell_1$ -penalty does not seem to improve results.

## 5 Discussion

We have shown that non-negative (or sign-constrained) least squares can be an effective regularization technique for sparse high-dimensional data under two conditions: (a) the data fulfil the so-called *Positive Eigenvalue Condition*, which is easy to check for a given dataset, and (b) the sign of the coefficients is known or can easily be estimated. If the conditions hold, NNLS can recover the correct sparsity pattern in the absence of any further shrinkage, as long as  $\log(p)s/n \rightarrow 0$  for  $n \rightarrow \infty$ , where  $p$  is the number of variables,  $s$  the number of non-zero variables in the optimal regression vector and  $n$  is sample size. We have shown network tomography as an example where the sign of regression coefficients is known a priori and the design condition is fulfilled automatically, at least approximately. In other examples the sign can be estimated by

an initial estimator. An attractive feature of NNLS is that it does not require any tuning parameter beyond the choice of the signs of the individual regression coefficients. Despite its simplicity, it can remarkably accurately for high-dimensional regression.

## 6 Appendix: Proofs

### 6.1 Proof of Theorem 1

First, for any  $C > 0$ ,  $1 - \Phi(C) \leq (2\pi)^{-1/2}C^{-1} \exp(-C^2/2)$ . Choosing  $C^2 = K_{p,\eta}^2 = 2 \log(\frac{\sqrt{2p}}{\sqrt{\pi\eta}})$ , it follows with  $\eta < 1/3$  and hence  $C \geq 1$  that  $1 - \Phi(C) \leq \eta/(2p)$ . Thus  $1 - (p+s)(1 - \Phi(C)) \geq 1 - (2p)(1 - \Phi(C)) \geq 1 - \eta$  and the results follow hence from Lemma 1.

### 6.2 Proof of Theorem 2

Define the oracle non-negative least squares solution as

$$\hat{\beta}^{oracle} := \operatorname{argmin}_{\beta} \|\mathbf{Y} - \mathbf{X}\beta\|_2^2 \quad \text{such that} \quad \min_k \beta_k \geq 0 \quad \text{and} \quad \beta_N \equiv 0, \quad (4)$$

and let  $\delta\beta = \hat{\beta} - \hat{\beta}^{oracle}$ .

Let  $M$  be the set  $M := \{k : \delta\beta_k < 0\}$ . Using Equation (9) in the proof of Lemma 1, it follows that, with probability at least  $1 - (p+s)(1 - \Phi(C))$ ,

$$\delta\beta^T \hat{\Sigma} \delta\beta \leq 2C\sigma \|\delta\beta_{M^c}\|_1 / \sqrt{n}$$

and, using  $\|\delta\beta_{M^c}\|_1 \leq \|\delta\beta\|_1$  and the bound in (7) for the latter quantity, it holds with probability at least  $1 - (p+s)(1 - \Phi(C))$  that

$$\delta\beta^T \hat{\Sigma} \delta\beta \leq 2C^2\sigma^2(5\nu^{-1} + 2\phi^{-1/2}) \frac{s}{n}.$$

Using again  $C^2 = K_{p,\eta}^2 = 2 \log(\frac{\sqrt{2p}}{\sqrt{\pi\eta}})$ , the claim follows.

### 6.3 Lemmata

**Lemma 1.** *Assume that the Positive Eigenvalue Condition holds with  $\nu > 0$ . Choose any  $C > 0$ . Assume that the compatibility condition holds with  $\phi > 0$  for  $L = 4\nu^{-1}$  and  $\min_{k \in S} \beta_k > C\sigma/\sqrt{n\phi}$ . It then holds with probability at least  $1 - (p+s)(1 - \Phi(C))$  that*

$$\|\hat{\beta} - \beta^*\|_1 \leq C\sigma(5/\nu + 4/\sqrt{\phi}) \frac{s}{\sqrt{n}}.$$

*Proof.* By the definition (1) of  $\hat{\beta}$  and definition (4) of  $\hat{\beta}^{oracle}$ ,

$$\delta\beta = \operatorname{argmin}_{\gamma} \|\mathbf{Y} - \mathbf{X}\hat{\beta}^{oracle} - \mathbf{X}\gamma\|_2^2 \quad \text{such that} \quad \gamma_k \geq -\hat{\beta}_k^{oracle} \quad \text{for all } k = 1, \dots, p. \quad (5)$$

The bound for  $\|\hat{\beta} - \beta^*\|_1$  follows as  $\|\hat{\beta} - \beta^*\|_1 \leq \|\hat{\beta}^{oracle} - \beta^*\|_1 + \|\delta\beta\|_1$ . Using Lemma 2, it holds with probability exceeding  $1 - (p+s)(1 - \Phi(C))$ ,

$$\|\hat{\beta}^{oracle} - \beta^*\|_1 \leq 2C\sigma\phi^{-1/2} \frac{s}{\sqrt{n}}, \quad (6)$$

and it thus remains to be shown that, if (6) is fulfilled, also

$$\|\delta\beta\|_1 \leq C\sigma(5\nu^{-1} + 2\phi^{-1/2}) \frac{s}{\sqrt{n}}. \quad (7)$$

Let  $\mathbf{R} = \mathbf{Y} - \mathbf{X}\hat{\beta}^{oracle}$ . Since  $\delta\beta \equiv 0$  is a feasible solution in (5), we have that

$$\delta\beta^T \mathbf{X}^T \mathbf{X} \delta\beta - 2\mathbf{R}^T \mathbf{X} \delta\beta \leq 0.$$

Let

$$M := \{k : \delta\beta_k < 0\} \quad (8)$$

By the definition of the estimator  $M \subseteq S$  and  $N \subseteq M^c$ . By Lemma 3, with probability at least  $1 - p(1 - \Phi(C))$ ,

$$\max_{k \in N} \mathbf{R}^T \mathbf{X}_k \leq C\sigma\sqrt{n}.$$

By Lemma 2, with probability at least  $1 - s(1 - \Phi(C))$ ,  $\mathbf{R}^T \mathbf{X}_k = 0$  for all  $k \in S$ . Hence, taken together, with probability at least  $1 - (p + s)(1 - \Phi(C))$ ,

$$\mathbf{R}^T \mathbf{X} \delta\beta \leq \left( \max_{k \in M^c} \mathbf{R}^T \mathbf{X}_k \right) \|\delta\beta_{M^c}\|_1 \leq C\sigma\sqrt{n} \|\delta\beta_{M^c}\|_1$$

and thus

$$\delta\beta^T \hat{\Sigma} \delta\beta \leq 2C\sigma \|\delta\beta_{M^c}\|_1 / \sqrt{n}. \quad (9)$$

Now,

$$\begin{aligned} \delta\beta^T \hat{\Sigma} \delta\beta &= \delta\beta_M^T \hat{\Sigma} \delta\beta_M + \delta\beta_{M^c}^T \hat{\Sigma} \delta\beta_{M^c} - 2 \sum_{i \in M, j \in M^c} \hat{\Sigma}_{i,j} \delta\beta_i \delta\beta_j \\ &\geq \delta\beta_{M^c}^T \hat{\Sigma} \delta\beta_{M^c} - 2 \|\delta\beta_M\|_1 \|\delta\beta_{M^c}\|_1 \\ &\geq \nu \|\delta\beta_{M^c}\|_1^2 - 2 \|\delta\beta_M\|_1 \|\delta\beta_{M^c}\|_1, \end{aligned} \quad (10)$$

having used the normalization to 1 of all columns of  $\mathbf{X}$  (which bounds the absolute values of all entries in  $\hat{\Sigma}$  by 1) for the second term in the second last inequality and the *Positive Eigenvalue Condition* for the first term in the last inequality (together with the fact that  $\min_{k \in M^c} \beta_k \geq 0$  by definition of  $M$  in (8)). Using this bound in (9) and dividing by  $\|\delta\beta_{M^c}\|_1$  yields that, with probability at least  $1 - (p + s)(1 - \Phi(C))$ ,

$$\begin{aligned} \|\delta\beta_{M^c}\|_1 &\leq 2\nu^{-1} \left( C\sigma/\sqrt{n} + \|\delta\beta_M\|_1 \right) \\ &= 2\nu^{-1} \left( \frac{C\sigma}{\|\delta\beta_M\|_1 \sqrt{n}} + 1 \right) \|\delta\beta_M\|_1 \end{aligned} \quad (11)$$

Evidently  $\|\delta\beta_M\|_1 \leq C\sigma/\sqrt{n}$  is either true or not. If it is true, then it follows trivially from the first inequality in (11) that  $\|\delta\beta_{M^c}\|_1 \leq 4\nu^{-1}C\sigma/\sqrt{n}$  and hence  $\|\delta\beta\|_1 \leq (1 + 4\nu^{-1})C\sigma/\sqrt{n} \leq 5\nu^{-1}C\sigma/\sqrt{n}$ , and the bound in (7) holds true.

Alternatively, if  $\|\delta\beta_M\|_1 > C\sigma/\sqrt{n}$  we have from the second inequality in (11) that  $\|\delta\beta_{M^c}\|_1 \leq L\|\delta\beta_M\|_1$  for  $L = 4\nu^{-1}$  and thus, using  $N \subseteq M^c$ , also  $\|\delta\beta_N\|_1 \leq L\|\delta\beta_S\|_1$ . The vector  $\delta\beta$  is then in  $\mathcal{R}(L, S)$ . Using the compatibility condition, it follows that

$$\delta\beta^T \hat{\Sigma} \delta\beta \geq \frac{\phi}{s} \|\delta\beta\|_1^2.$$

Using this in (9),

$$\frac{\phi}{s} \|\delta\beta\|_1^2 \leq 2C\sigma \|\delta\beta_{M^c}\|_1 / \sqrt{n}. \quad (12)$$

Using  $\|\delta\beta_{M^c}\|_1 \leq \|\delta\beta\|_1$ , it follows that

$$\|\delta\beta\|_1 \leq 2C\sigma s / \sqrt{n\phi}, \quad (13)$$

which also satisfies the bound in (7). Hence, the bound (7) holds under both possible scenarios ( $\|\delta\beta_M\|_1 \leq C\sigma/\sqrt{n}$  true or false) and the proof is complete.  $\square$

**Lemma 2.** Let  $\hat{\beta}^{ols}$  be the least squares estimator restricted to  $S$ :

$$\hat{\beta}^{ols} = \operatorname{argmin}_{\beta} \|\mathbf{Y} - \mathbf{X}\beta\|_2^2 \text{ such that } \beta_N \equiv 0.$$

If  $\min_{k \in S} \beta_k^* \geq C\sigma/\sqrt{n\phi}$ , Then

$$P(\hat{\beta}^{ols} \equiv \hat{\beta}^{oracle}) \geq 1 - s(1 - \Phi(C)),$$

and, with at least the same probability  $1 - s(1 - \Phi(C))$ ,

$$\|\beta^* - \hat{\beta}^{oracle}\|_{\infty} \leq C\sigma/\sqrt{n\phi}$$

*Proof.* It is only necessary to show that  $\min_{k \in S} \hat{\beta}_k^{ols} \geq 0$  with probability at least  $1 - s(1 - \Phi(C))$ .

The error term has, under the made assumptions, a normal distribution,  $\hat{\beta}_k^{ols} - \beta_k^* \sim \mathcal{N}(0, \sigma^2(n\hat{\Sigma}_{SS})_k^{-1})$  for all  $k \in S$ . The minimal eigenvalue of  $\hat{\Sigma}_{SS}$  is bounded from below by  $\phi$  by the compatibility condition and the variance of  $\hat{\beta}_k^{ols}$  is thus bounded from above by  $\phi^{-1}\sigma^2/n$  for all  $k \in S$ . It follows with Bonferroni's inequality that, with probability at least  $1 - s(1 - \Phi(C))$ ,

$$\|\beta^* - \hat{\beta}^{ols}\|_{\infty} \leq C\sigma/\sqrt{n\phi}. \quad (14)$$

If  $\min_{k \in S} \beta_k^* \geq C\sigma/\sqrt{n\phi}$ , then (14) implies that  $\min_{k \in S} \hat{\beta}_k^{ols} \geq 0$  and thus  $\hat{\beta}^{oracle} \equiv \hat{\beta}^{ols}$  and thus also

$$\|\beta^* - \hat{\beta}^{oracle}\|_{\infty} \leq C\sigma/\sqrt{n\phi},$$

which completes the proof.  $\square$

**Lemma 3.** With probability at least  $1 - p(1 - \Phi(C))$ ,

$$\max_{k \in N} (\mathbf{Y} - \mathbf{X}\hat{\beta}^{oracle})^T \mathbf{X}_k \leq C\sigma\sqrt{n}.$$

*Proof.* We condition on the event  $\hat{\beta}^{oracle} \equiv \hat{\beta}^{ols}$ , which happens according to Lemma 2 with probability at least  $1 - s(1 - \Phi(C))$ . Then  $\mathbf{Y} - \mathbf{X}\hat{\beta}^{oracle} = \mathbf{Y} - \mathbf{X}\hat{\beta}^{ols} = P_{S^{\perp}}\mathbf{Y}$ , where  $P_{S^{\perp}}\mathbf{Z}$  is the projection of a vector  $\mathbf{Z} \in \mathbb{R}^n$  into the space orthogonal to  $\mathbf{X}_S$ . Now,  $P_{S^{\perp}}\mathbf{Y} = P_{S^{\perp}}(\mathbf{X}\beta^* + \varepsilon) = P_{S^{\perp}}\varepsilon$ . The distribution of  $(P_{S^{\perp}}\varepsilon)^T \mathbf{X}_k$  is, for every  $k \in N$ , normal with mean 0 and variance at most  $\sigma^2 n$ , and thus  $P((P_{S^{\perp}}\varepsilon)^T \mathbf{X}_k \geq C\sigma\sqrt{n}) \leq 1 - \Phi(C)$  for all  $k \in N$  and, using a Bonferroni bound,  $P(\max_{k \in N} (P_{S^{\perp}}\varepsilon)^T \mathbf{X}_k \geq C\sigma\sqrt{n}) \leq |N|(1 - \Phi(C))$ . The unconditional probability of  $\max_{k \in N} (P_{S^{\perp}}\varepsilon)^T \mathbf{X}_k \geq C\sigma\sqrt{n}$  is thus at least  $1 - s(1 - \Phi(C)) - |N|(1 - \Phi(C)) = 1 - (s + |N|)(1 - \Phi(C)) = 1 - p(1 - \Phi(C))$ , which completes the proof.  $\square$

## References

- S. Bellavia, M. Macconi, and B. Morini. An interior point newton-like method for non-negative least-squares problems with degenerate solution. *Numerical Linear Algebra with Applications*, 13:825–846, 2006.
- P. Bickel, Y. Ritov, and A. Tsybakov. Simultaneous analysis of Lasso and Dantzig selector. *Annals of Statistics*, 37:1705–1732, 2009.
- C. Boutsidis and P. Drineas. Random projections for the nonnegative least-squares problem. *Linear Algebra and its Applications*, 431:760–771, 2009.
- L. Breiman. Better subset regression using the nonnegative garrote. *Technometrics*, 37:373–384, 1995.
- B.F. Bunea, A.B. Tsybakov, and M.H. Wegkamp. Aggregation for gaussian regression. *Annals of Statistics*, 35:1674–1697, 2007.

- F. Bunea, A. Tsybakov, and M. Wegkamp. Sparsity oracle inequalities for the lasso. *Electronic Journal of Statistics*, pages 169–194, 2006.
- E. Candes and T. Tao. The Dantzig selector: statistical estimation when  $p$  is much larger than  $n$ . *Annals of Statistics*, 35:2312–2351, 2007.
- R. Castro, M. Coates, G. Liang, R. Nowak, and B. Yu. Network tomography: Recent developments. *Statistical Science*, pages 499–517, 2004.
- D. Chen and R.J. Plemmons. *Nonnegativity constraints in numerical analysis*. World Scientific Press, River Edge, NJ, USA, 2009.
- C.H.Q. Ding, T. Li, and M.I. Jordan. Convex and semi-nonnegative matrix factorizations. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 32:45–55, 2010.
- D.L. Donoho, I.M. Johnstone, J.C. Hoch, and A.S. Stern. Maximum entropy and the nearly black object. *Journal of the Royal Statistical Society. Series B*, pages 41–81, 1992.
- B. Efron, T. Hastie, I. Johnstone, and R. Tibshirani. Least angle regression. *Annals of Statistics*, 32:407–451, 2004.
- A.E. Hoerl and R.W. Kennard. Ridge regression: Biased estimation for nonorthogonal problems. *Technometrics*, pages 55–67, 1970.
- D. Kim, S. Sra, and I.S. Dhillon. A new projected quasi-newton approach for the nonnegative least squares problem. Technical report, Department of Computer Science, University of Texas, TR-06-54, 2006.
- H. Kim and H. Park. Sparse non-negative matrix factorizations via alternating non-negativity-constrained least squares for microarray data analysis. *Bioinformatics*, 23:1495, 2007.
- E. Lawrence, G. Michailidis, and V.N. Nair. Network delay tomography using flexicast experiments. *Journal of the Royal Statistical Society: Series B*, 68:785–813, 2006.
- C.L. Lawson and R.J. Hanson. *Solving least squares problems*, volume 15. Society for Industrial Mathematics, 1995.
- D.D. Lee and H.S. Seung. Algorithms for non-negative matrix factorization. *Advances in neural information processing systems*, 13, 2001.
- D.D. Lee, H.S. Seung, et al. Learning the parts of objects by non-negative matrix factorization. *Nature*, 401:788–791, 1999.
- N. Meinshausen and B. Yu. Lasso-type recovery of sparse representations from high-dimensional data. *Annals of Statistics*, 7:246–270, 2009.
- M. Slawski, M. Hein, and E. Campus. Sparse recovery by thresholded non-negative least squares. Technical report, Department of Computer Science, University of Saarbruecken, 2011.
- R. Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society, Series B*, 58:267–288, 1996.
- S.A. Van De Geer. High-dimensional generalized linear models and the lasso. *The Annals of Statistics*, 36:614–645, 2008.
- S.A. Van De Geer and P. Bühlmann. On the conditions used to prove oracle results for the lasso. *Electronic Journal of Statistics*, 3:1360–1392, 2009.
- M.J. Wainwright. Sharp thresholds for high-dimensional and noisy recovery of sparsity. *IEEE Transactions on Information Theory*, 55:2183–2202, 2009.

- M. S. Waterman. Least squares with nonnegative regression coefficients. *Journal of Statistical Computation and Simulation*, 6:67–70, 1977.
- B. Xi, G. Michailidis, and V.N. Nair. Estimating network loss rates using active tomography. *Journal of the American Statistical Association*, 101:1430–1448, 2006.
- M. Yuan and Y. Lin. Model selection and estimation in regression with grouped variables. *Journal of the Royal Statistical Society: Series B*, 68:49–67, 2006.
- C.H. Zhang and J. Huang. The sparsity and bias of the lasso selection in high-dimensional linear regression. *Annals of Statistics*, 36:1567–1594, 2008.
- H. Zou. The adaptive lasso and its oracle properties. *Journal of the American Statistical Association*, 101:1418–1429, 2006.