

# Improving Lasso for Model Selection and Prediction

**Piotr Pokarowski**

POKAR@MIMUW.EDU.PL

**Agnieszka Prochenka**

A.PROCHENKA@PHD.IPIPAN.WAW.PL

**Micha Frej**

M.FREJ@MIMUW.EDU.PL

**Wojciech Rejchel**

WREJCHEL@GMAIL.COM

*Faculty of Mathematics, Informatics and Mechanics, University of Warsaw, Banacha 2, 02-097 Warsaw, Poland*

**Jan Mielniczuk**

MIEL@IPIPAN.WAW.PL

*Institute of Computer Sciences, Jana Kazimierza 5, 01-248 Warsaw, Poland*

## Abstract

The Lasso, that is  $l_1$ -penalized loss estimator is a popular tool for fitting sparse models to high-dimensional data. The concave regularizations SCAD or MCP approximate more closely  $l_0$ -penalized loss, that is the Generalized Information Criterion (GIC), and correct intrinsic estimation bias of the Lasso. In this paper we propose an alternative method of improving the Lasso for predictive models which encompass linear and logistic models as premier examples. The approach, for a given penalty, orders the absolute values of the Lasso non-zero coefficients and then selects the model from a small nested family by GIC. We derive an upper bound on the methods selection error and show in numerical experiments on synthetic and real-world data sets that an implementation of our algorithm is more accurate than implementations of studied concave regularizations.

**Keywords:** List of keywords

## 1. Introduction

Sparse high-dimensional predictive models, where the number of true predictors  $t$  is significantly smaller than the sample size  $n$  and the number of all predictors  $p$  greatly exceeds  $n$  have been a focus of research in statistical machine learning in recent years. The Lasso algorithm, that is the minimum loss method regularized by sparsity inducing  $l_1$  penalty, is the main tool of fitting such models (Tibshirani, 2011; Bühlmann and van de Geer, 2011). However, a few years ago it has been shown that the model selected by the Lasso is usually too large and that for asymptotically consistent model selection it requires the *irrepresentable condition* on an experimental matrix (Meinshausen and Bühlmann, 2006; Zhao and Yu, 2006; Shen et al., 2012) which is too restrictive in general. Model's dimension can be reduced without loss of quality using the Thresholded Lasso (TL) algorithm, which selects variables with largest absolute values of the Lasso coefficients (Ye and Zhang, 2010; Zhou, 2009) or by solving a more computationally demanding minimization of a loss with a folded concave penalty (FCP) as SCAD (Fan and Li, 2001), MCP (Zhang, 2010a) or capped  $l_1$ -penalty (Zhang, 2010b; Shen et al., 2012). TL, FCP and similar methods lead to consistent selection under weaker assumptions such as the *restricted isometry property* (Zhang, 2010a,b; Zhang and Zhang, 2012; Shen et al., 2012; Wang et al., 2013; Fan et al., 2014; Wang et al., 2014). Recently in Pokarowski and Mielniczuk (2015) an algorithm called *Screening–Ordering–Selection*

(SOS) was introduced for linear model selection, which reduces the model selected by the Lasso. SOS is based on the variant of TL proposed by [Zhou \(2009\)](#).

The first main contribution of the paper is that we simplify and generalize SOS in the two step *Screening–Selection* (SS) algorithm for general predictive models e.g. linear normal, logistic or Cox proportional hazard models. In the first screening step, one computes the Lasso estimator  $\hat{\beta}$  with penalty  $\lambda$  and orders its nonzero coefficients according to their decreasing absolute values. In the second, selection step, one chooses the model which minimizes GIC with penalty  $\lambda^2/2$  in a nested family induced by the ordering. Post-model selection estimator of  $\beta$  is the minimum loss estimator for the chosen model. Thus the SS algorithm ([Algorithm 1](#) below) is the Lasso with adaptive thresholding based on GIC. We derive an exponential upper bound on selection error of SS in terms of  $\lambda$  ([Theorem 2](#)), which parallels the known bounds for TL, see [Theorem 8](#) for linear models in [Ye and Zhang \(2010\)](#) or FCP, see [Corollary 3 and 5](#) in [Fan et al. \(2014\)](#). However, in contrast to these methods, SS is constructive for linear models in that it relies neither on the unknown parameters as the true vector  $\beta$  or the cone invertibility factors. Instead,  $\lambda$  only depends on the sample size, the number of predictors and an upper bound on the noise parameter.

Although TL, FCP, SOS or SS algorithms use the Lasso estimators only for one value of the penalty, which is convenient for theoretical analysis, the practical Lasso implementations return coefficient estimators for all possible penalty values (algorithm and R package `LARS` described in [Efron et al., 2004](#)) or for a given net of them (R package `glmnet` described in [Friedman et al. \(2010\)](#)). Similarly, using a net of penalty values, the FCP algorithm has been implemented for linear (in R package `SparseNet` [Mazumder et al. \(2011\)](#)) and logistic models (in R package `cvplogistic` [Jiang and Huang \(2014\)](#)).

Our second main contribution is that we propose the SOSnet algorithm ([Algorithm 2](#) below), which is a generalization of the SOS algorithm for general predictive models. SOSnet uses `glmnet` for a net of penalty values, then for each of them it orders the chosen predictors according to Wald statistics and finally, it selects the model from a small family by minimizing GIC. We show in numerical experiments for linear and logistic models on synthetic and real-world data sets that SOSnet is more accurate than implementations of FCP.

## 2. Models and Fitting Algorithms

In this section we start with definitions of considered models and estimation criteria, next we present model selection algorithms.

### 2.1. Models

The way we model data will encompass normal linear and logistic models as premier examples. Our assumptions are stated in their most general form which allows proving exponential bounds for probability of selection error without obscuring their essentiality. We consider independent data  $(y_1, x_1), (y_2, x_2), \dots, (y_n, x_n)$ , where  $y_i \in \mathbb{R}$ ,  $x_i \in \mathbb{R}^p$  for  $i = 1, 2, \dots, n$ , a known differentiable *cumulant* function  $\gamma : \mathbb{R} \rightarrow \mathbb{R}$  and we assume that for some *true*  $\beta \in \mathbb{R}^p$

$$\mathbb{E}y_i = \dot{\gamma}(x_i^T \beta) \text{ for } i = 1, 2, \dots, n. \quad (1)$$

Note that (1) is satisfied in particular by the Generalized Linear Models (GLM) and a nonlinear regression with an additive error. We will use also vectorized version of the cumulant and its

derivative. For  $\eta = (\eta_1, \dots, \eta_n)^T$  we define  $\gamma(\eta) = (\gamma(\eta_1), \dots, \gamma(\eta_n))^T$  and similarly  $\dot{\gamma}(\eta) = (\dot{\gamma}(\eta_1), \dots, \dot{\gamma}(\eta_n))^T$ .

Let  $X = [x_{\cdot 1}, \dots, x_{\cdot p}] = [x_1, \dots, x_n]^T$  be a  $n \times p$  matrix of experiment and  $J \subseteq \{1, 2, \dots, p\} = F$  be an arbitrary subset of the *full model*  $F$ ,  $\bar{J} = F \setminus J$ . As  $J$  may be viewed as sequence of zeros and ones on  $F$ ,  $|J| = |J|_1$  denotes cardinality of  $J$ . Let  $\beta_J$  be a subvector of  $\beta$  with elements having indices in  $J$ ,  $X_J$  be a submatrix of  $X$  with columns having indices in  $J$  and  $r(X_J)$  denotes a rank of  $X_J$ . Linear model pertaining to predictors being columns of  $X_J$  will be frequently identified as  $J$ . In particular, let  $T$  denotes a *true model* that is  $T = \text{supp}(\dot{\beta}) = \{j \in F : \dot{\beta}_j \neq 0\}$  and  $t = |T|$ .

We assume also that a *total cumulant* function  $g(\beta) = \sum_{i=1}^n \gamma(x_i^T \beta)$  is convex and, additionally, *strongly convex* at  $\dot{\beta}$  in a sense that exists  $c \in (0, 1]$  such that for all *sparse*  $\beta \in \mathbb{B} \equiv \mathbb{B}(X, \dot{\beta}, \bar{t})$  we have

$$g(\beta) \geq g(\dot{\beta}) + (\beta - \dot{\beta})^T \dot{g}(\dot{\beta}) + \frac{c}{2} (\dot{\beta} - \beta)^T X^T X (\dot{\beta} - \beta), \quad (2)$$

where  $t \leq \bar{t} < n \wedge p$ ,

$$\mathbb{B} = \bigcup_{J \supset T, r(X_J) = |J| \leq \bar{t}} \{\beta_J : \|X_T(\dot{\beta}_J - \beta_J)\|^2 \leq \delta_{t-1}\}, \quad (3)$$

$$\delta_{t-1} = \min_{j \in T} \|X_T \dot{\beta}_T - x_{\cdot j} \dot{\beta}_j\|^2. \quad (4)$$

We note that this crucial property of the total cumulant is slightly weaker than an usual definition of strong convexity which would have a second derivative of  $g$  at  $\dot{\beta}$  in place of  $X^T X$ . Let us remark that  $\dot{g}(\beta) = X^T \dot{\gamma}(X\beta)$ .

Moreover, we assume that centred responses  $\varepsilon_i = y_i - \mathbb{E}y_i$  have a *subgaussian distributions* with the same constant  $\sigma$ , that is for  $i = 1, 2, \dots, n$  and  $u \in \mathbb{R}$  we have

$$\mathbb{E} \exp(u\varepsilon_i) \leq \exp(\sigma^2 u^2 / 2). \quad (5)$$

**Examples.** For a normal linear model  $y_i \sim N(x_i^T \dot{\beta}, \dot{\sigma}^2)$ ,  $\gamma(\eta_i) = \eta_i^2 / 2$  and (5) is fulfilled with any  $\sigma \geq \dot{\sigma}$ . For a logistic model  $y_i \sim \text{binom}(1, [1 + \exp(-x_i^T \dot{\beta})]^{-1})$ ,  $\gamma(\eta_i) = \log(1 + \exp(\eta_i))$  and as  $(\varepsilon_i)$  are bounded random variables, then (5) is satisfied with any  $\sigma \geq 1/2$ . It is easy to note that for a linear model the strong convexity assumption (2) is fulfilled with  $c = 1$  and for a logistic model with

$$c = \min_i \min_{\beta \in \mathbb{B}} \exp(x_i^T \beta) / (1 + \exp(x_i^T \beta))^2.$$

## 2.2. Fitting Algorithms

For estimation of  $\dot{\beta}$  we consider a *loss function*

$$\ell(\beta) = \sum_{i=1}^n (\gamma(x_i^T \beta) - y_i x_i^T \beta) = g(\beta) - \beta^T X^T y, \quad (6)$$

where  $y = (y_1, \dots, y_n)^T$ . It is easy to see that  $\dot{\ell}(\beta) = X^T (\dot{\gamma}(X\beta) - y)$ , and consequently  $\dot{\ell}(\dot{\beta}) = -X^T \varepsilon$  for  $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)^T$ . Moreover, observe that  $\dot{\beta} = \text{argmin}_{\beta} \mathbb{E} \ell(\beta)$  and (2) is equivalent to strong convexity of  $\ell$  in  $\dot{\beta}$  for all  $\beta \in \mathbb{B}$

$$\ell(\beta) \geq \ell(\dot{\beta}) + (\dot{\beta} - \beta)^T X^T \varepsilon + \frac{c}{2} (\dot{\beta} - \beta)^T X^T X (\dot{\beta} - \beta). \quad (7)$$

---

**Algorithm 1** TL and SS
 

---

**Input:**  $y, X$  and  $\lambda, \tau$ 
**Screening (Lasso)**  $\hat{\beta} = \operatorname{argmin}_{\beta} \{\ell(\beta) + \lambda|\beta|_1\}$ ;

**Thresholded Lasso**  $\hat{T}_{TL} = \{j : |\hat{\beta}_j| > \tau\}$ ;

**Selection (GIC)**

 order nonzero  $|\hat{\beta}_{j_1}| \geq \dots \geq |\hat{\beta}_{j_s}|$ , where  $s = |\operatorname{supp}\hat{\beta}|$ ;

 set  $\mathcal{J} = \{\{j_1\}, \{j_1, j_2\}, \dots, \operatorname{supp}\hat{\beta}\}$ ;

 $\hat{T}_{SS} = \operatorname{argmin}_{J \in \mathcal{J}} \{\ell_J + \lambda^2/2|J|\}$ .

**Output:**  $\hat{T}_{TL}, \hat{T}_{SS}$ 


---

Let  $\hat{\beta}_J^{ML} = \operatorname{argmin}_{\beta_J} \ell(\beta_J)$  denotes a *minimum loss estimator* based on  $y$  and  $\{x_{\cdot j}, j \in J\}$  and  $\ell_J = \ell(\hat{\beta}_J^{ML})$ . Note that for GLM  $\hat{\beta}_J^{ML}$  coincides with a maximum likelihood estimator for a model pertaining to  $J$ . Finally, for  $\beta \in \mathbb{R}^p$  and  $q \geq 1$  let  $|\beta|_q = (\sum_{j=1}^p |\beta_j|^q)^{1/q}$  be  $\ell_q$  norm.

In Algorithm 1 we present two selection procedures: the first one is the Thresholded Lasso (TL) method which consists of retaining only these variables for which absolute values of their Lasso estimators exceed a certain threshold  $\tau$ . The second one, named Screening–Selection (SS) procedure finds minimal value of Generalized Information Criterion (GIC) for the nested family which is constructed using ordering of the nonzero Lasso estimates.

### 3. Selection Error Bounds for TL and SS

In this section we present upper exponential bounds on the selection error of the TL and SS algorithms. In order to make the exposition simpler we assume that columns of  $X$  are normalized in such a way that  $\|x_j\| = 1$  for  $j = 1, \dots, p$ . Moreover, let  $\hat{\beta}_{min} = \min_{j \in T} |\hat{\beta}_j|$ .

#### 3.1. A Bound for TL

First we generalize a characteristic of linear models which quantifies the degree of separation between the true model  $T$  and other models introduced in [Ye and Zhang \(2010\)](#). For  $a \in (0, 1)$  consider a signed pseudo-cone

$$\mathcal{C}_a = \left\{ \nu \in \mathbb{R}^p : |\nu_{\bar{T}}|_1 \leq \frac{1+a}{1-a} |\nu_T|_1, \nu_j x_j^T [\dot{\gamma}(X(\hat{\beta} + \nu)) - \dot{\gamma}(X\hat{\beta})] \leq 0, j \in \bar{T} \right\}. \quad (8)$$

For  $q \geq 1$  and  $a \in (0, 1)$  let a Sign-Restricted Pseudo-Cone Invertibility Factor (SCIF) be defined as

$$\zeta_{a,q} = \inf_{\nu \in \mathcal{C}_a} \frac{|X^T [\dot{\gamma}(X(\hat{\beta} + \nu)) - \dot{\gamma}(X\hat{\beta})]|_{\infty}}{|\nu|_q} \quad (9)$$

We let  $\zeta_a = \zeta_{a,\infty}$ . In comparison to more popular restricted eigenvalues or compatibility constants, variants of SCIF enable sharper  $\ell_q$  estimation error bounds of the Lasso for  $q > 2$  (cf Corollary 4 in [Ye and Zhang \(2010\)](#) and [Huang and Zhang \(2012\)](#); [Zhang and Zhang \(2012\)](#)).

The following lemma is a generalization of Theorem 3 in [Ye and Zhang \(2010\)](#).

**Lemma 1** *If  $\ell$  is convex and  $a \in (0, 1)$ , then on  $\{|X^T \varepsilon|_\infty \leq a\lambda\}$  we have  $|\hat{\beta} - \mathring{\beta}|_q \leq (1+a)\lambda\zeta_{a,q}^{-1}$ .*

**Theorem 2** *If  $\ell$  is convex,  $(\varepsilon_i)_i$  are subgaussian with  $\sigma$  and for  $a_1, a_2 \in (0, 1)$  we have*

$$2a_1^{-2}a_2^{-1}\sigma^2 \log p \leq \lambda^2 \leq (1+a_1)^{-2}\zeta_{a_1}^2 \tau^2 < (1+a_1)^{-2}\zeta_{a_1}^2 \mathring{\beta}_{min}^2/4,$$

then

$$\mathbb{P}(\hat{T}_{TL} \neq T) \leq 2 \exp\left(-\frac{(1-a_2)a_1^2\lambda^2}{2\sigma^2}\right). \quad (10)$$

Constant  $a_2$  is used to remove multiplicative factor  $p$  from the exponential bound at the expense of slightly diminishing the exponent in (10). Note that assumptions of Theorem 2 stipulate that truncation level  $\tau$  is contained in interval  $[(1+a_1)\lambda\zeta_{a_1}^{-1}, \mathring{\beta}_{min}/2)$ . Analogous theorem for FCP, see Corollary 3 for linear models and Corollary 5 for logistic models in Fan et al. (2014), requires an additional assumption on the minimal eigenvalue of  $X_T^T X_T$  and the proof is more difficult, but for both methods a condition on  $\tau$  requires unknown  $\zeta_a$  or  $\mathring{\beta}_{min}$ .

### 3.2. A Bound for SS

Let  $H_J$  be an orthogonal projection matrix onto the subspace spanned by columns of  $X_J$ . A scaled K-L distance between  $T$  and its submodels, see Shen et al. (2012, 2013) is

$$\delta = \min_{J \subset T} \frac{\|(I - H_J)X_T \mathring{\beta}_T\|^2}{|T \setminus J|}. \quad (11)$$

Different variants of the K–L distance have been often used in the consistency analysis of selection algorithms, cf Section 3.1 in Pokarowski and Mielniczuk (2015), but  $\delta$  defined above seems to lead to optimal results, cf Theorem 1 in Shen et al. (2013). Technical constants  $a_1, \dots, a_4$  allow to avoid *ad-hoc* coefficients in the bounds and simplify asymptotic considerations. For given  $1/2 < a_1 < 1$  define  $a_2 = 1 - (1 - \log(1 - a_1))(1 - a_1)$ ,  $a_3 = 2 - 1/a_1$  and  $a_4 = \sqrt{a_1 a_2}$ . Note that  $a_2, a_3$  and  $a_4$  are functions of  $a_1$  and obviously if  $a_1 \rightarrow 1$ , then  $a_2, a_3, a_4 \rightarrow 1$ .

**Theorem 3** *Assume (1)–(5) and that for  $a_1 \in (1/2, 1)$*

$$\frac{2\sigma^2 \log p}{a_3 a_2 a_1 c} \vee \frac{\sigma^2 t}{(1-a_1)^2 c} \leq \lambda^2 < \frac{\delta c}{(1 + \sqrt{2(1-a_1)})^2} \wedge \frac{\zeta_{a_4}^2 \mathring{\beta}_{min}^2}{4(1+a_4)^2}. \quad (12)$$

Then

$$\mathbb{P}(\hat{T}_{SS} \neq T) \leq 4.5 \exp\left(-\frac{a_2(1-a_1)c\lambda^2}{2\sigma^2}\right). \quad (13)$$

Selection consistency that is asymptotic correctness of  $\hat{T}_{SS}$  now easily follows.

**Corollary 4** *Assume that  $t = o(\log p)$  for  $n \rightarrow \infty$  and set  $a_1 = 1 - \sqrt{\frac{t}{2 \log p}}$ ,  $\lambda^2 = \frac{2\sigma^2 \log p}{a_3 a_2 a_1 c}$ . Then  $\lambda^2 = c^{-1}(2\sigma^2 \log p)(1 + o(1))$ . If additionally  $\mathring{\beta}$  is asymptotically identifiable:*

$$\overline{\lim}_n \frac{2\sigma^2 \log p}{(c^2 \delta) \wedge (\zeta_{a_4}^2 \mathring{\beta}_{min}^2 / 16)} < 1, \quad \text{then } \mathbb{P}(\hat{T}_{SS} \neq T) = o(1).$$

**Remarks. 1.** Theorem 3 for linear models may be obtained directly, analogously as in Pokarowski and Mielniczuk (2015). The resulting lower bound on  $\lambda^2$  is  $\lambda^2 \geq 2\sigma^2 \log p / (a_1 a_2 a_3)$  without the additional condition  $\lambda^2 \geq \sigma^2 t / (1 - a_1)^2$  assumed in Theorem 3 (recall that  $c = 1$ ).

**2.** Theorem 2 determines the interval of admissible  $\lambda$ , a parameter of SS, for which a bound in (13) holds. Corollary 1 states more easily interpretable result: for  $\lambda$  equal to lower endpoint of the above interval SS is asymptotically correct provided that the true model is asymptotically identifiable. Although identifiability condition is not effectively verifiable,  $\lambda$  can be explicitly given for linear models as

$$\lambda = \sqrt{2\sigma^2 \log p} (1 + o(1)) \tag{14}$$

and for logistic models as

$$\lambda = \sqrt{(\log p) / (2c)} (1 + o(1)), \tag{15}$$

since  $\sigma \geq 1/2$ . Thus for linear models a parameter of SS is given constructively in contrast to TL or FCP which require an additional parameter  $\tau$ , depending on identifiability constants as SCIF. In literature concerning Lasso and its modifications the smallest possible  $\lambda$  is taken as the default value as then the algorithm is asymptotically consistent for the largest class of models (the same approach is adopted for prediction and estimation). Such  $\lambda$  will be called the *safest choice* in the Conclusions.

#### 4. Extension to general convex contrasts

In this part of the paper we investigate properties of the SS algorithm beyond GLM. In fact, the main assumption, that will be required, is convexity of the "contrast function". We show that the SS algorithm is very flexible procedure that can be (succesfully) applied to the various spectrum of practical problems.

First, for  $\beta \in \mathbb{R}^p$  and a contrast function  $\phi : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$  we define a loss function

$$\ell(\beta) = \sum_{i=1}^n \phi(\beta^T x_i, y_i).$$

Considering the standard linear model one usually uses the quadratic contrast  $\phi(\beta^T x_i, y_i) = (y_i - \beta^T x_i)^2$  as we have done before. However, it is well known that the quadratic contrast is very sensitive to the distribution of errors  $\varepsilon_i$  and does not work well, if this distribution is, for instance, heavy-tailed and outliers appear. To overcome this difficulty we can use the absolute contrast  $\phi(\beta^T x_i, y_i) = |y_i - \beta^T x_i|$ . Next, working with dichotomous  $y_i$  we can apply the logistic regression that belongs to GLM and has been considered ealier. In this case we have  $\phi(\beta^T x_i, y_i) = -y_i \beta^T x_i + \log[1 + \exp(\beta^T x_i)]$ . But there are also very popular and efficient algorithms called support vector machines that use, for instance, the following contrast  $\phi(\beta^T x_i, y_i) = [\max(0, 1 - y_i \beta^T x_i)]^2$ .

Our main assumption is that the contrast function  $\phi$  is convex with respect to  $\beta$ . All examples in the above paragraph satisfy this property. Notice that they need not be differentiable nor decompose, as in (6) for GLM, into the sum of the nonrandom cumulant  $\gamma$  and the random linear term  $y_i \beta^T x_i$ . In this section we prove the exponential upper bound of the error of the SS algorithm based on these contrast functions that is almost the same as in Theorem 3. To do it we strongly exploit the empirical process theory. Obviously, the SS algorithm is the same as in Algorithm 1.

We add few definitions and notations to those in the previous parts of the paper. We start with defining two balls: the first one is the  $l_1$ -ball  $B_1(r) = \{\beta : |\beta - \hat{\beta}|_1 \leq r\}$  with radius  $r > 0$ . The

second one is the  $l_2$ -ball  $B_{2,J}(r) = \{\beta_J : \|X_J(\hat{\beta} - \beta_J)\|^2 \leq r^2\}$  with radius  $r > 0$ , where  $J$  is a (sparse) subset of  $\{1, \dots, p\}$  and  $T \subset J$ . Recall that  $\hat{\beta}$  is, as previously, a minimizer  $\mathbb{E}\ell(\beta)$ . Besides, let  $B_J = B_J(\sqrt{\delta_{t-1}})$ , where  $\delta_{t-1}$  is defined in (4). In further argumentation key roles are played by:

$$Z(r) = \sup_{\beta \in B_1(r)} \left| \ell(\beta) - \mathbb{E}\ell(\beta) - [\ell(\hat{\beta}) - \mathbb{E}\ell(\hat{\beta})] \right|$$

and

$$U_J(r) = \sup_{\beta \in B_{2,J}(r)} \left| \ell(\beta) - \mathbb{E}\ell(\beta) - [\ell(\hat{\beta}) - \mathbb{E}\ell(\hat{\beta})] \right|,$$

which are empirical processes over  $l_1$  and  $l_2$ -balls, correspondingly. We need also the compatibility factor that is an analog of SCIF defined in (9). Namely, for arbitrary  $a \in (0, 1)$  a compatibility factor is

$$\kappa_a = \inf_{0 \neq \beta \in \mathcal{C}_a} \frac{\beta^T X^T X \beta}{|\beta_T|_1^2}, \quad (16)$$

where  $\mathcal{C}_a$  is a simplified version of (8), namely

$$\mathcal{C}_a = \left\{ \nu \in \mathbb{R}^p : |\nu_T|_1 \leq \frac{1+a}{1-a} |\nu_T|_1 \right\}.$$

Convexity of the contrast function is the main assumption in this section. However, similarly to the previous section we need also the following *strong convexity* of  $\mathbb{E}\ell(\beta)$  at  $\hat{\beta}$ : there exists  $c_1 \in (0, 1]$  ( $c_2 \in (0, 1]$ , respectively) such that for each  $\beta_1 \in B_1(\hat{\beta}_{min})$  ( $\beta_2 \in B$ , respectively) we have for  $i = 1, 2$

$$\mathbb{E}\ell(\beta_i) - \mathbb{E}\ell(\hat{\beta}) \geq \frac{c_i}{2} (\beta_i - \hat{\beta})^T X^T X (\beta_i - \hat{\beta}). \quad (17)$$

Notice that we require the expected loss  $\mathbb{E}\ell(\beta)$ , not the loss  $\ell(\beta)$ , to be strongly convex. Besides, while considering GLM the condition (7) is equivalent to (17) for  $i = 2$ . Finally, to prove exponential bounds for GLM we use subgaussianity that allows us to obtain probabilistic inequalities in Lemma 10. In this section we need the analog of (23) of the form: there exists  $L > 0$  and constants  $K_1, K_2 > 0$  such that for each  $0 < r \leq \hat{\beta}_{min}$  and  $z \geq 1$  we have

$$P\left(\frac{Z(r)}{r} > K_1 L z \sqrt{\log(2p)}\right) \leq \exp(-K_2 \log(2p) z^2). \quad (18)$$

Besides, the inequality (24) is replaced by the following: there exists  $L > 0$  and constants  $K_1, K_2 > 0$  such that for each  $0 < r \leq \sqrt{\delta_{t-1}}$ ,  $z \geq 1$  and  $J : T \subset J, r(X_J) = |J| \leq \bar{t}$  we have

$$P\left(\frac{U_J(r)}{r} > K_1 L z \sqrt{|J|}\right) \leq \exp(-K_2 |J| z^2). \quad (19)$$

The detailed comparison between assumptions and results for models in this section and those for GLM is given after the main result of this section, which is now stated.

**Theorem 5** Fix  $a \in (0, 1)$ . Assume that (17), (18), (19) and

$$K_1 \max(\log p, t/a) \frac{L^2}{c_2} \leq \lambda^2 \leq K_2 \min \left[ \frac{c_2 \delta}{(1 + \sqrt{2a})^2}, \frac{c_2 \delta_{t-1}}{\bar{t} - t}, \left( c_1 \kappa_a \hat{\beta}_{min} \right)^2 \right]. \quad (20)$$

Then

$$P\left(\hat{T}_{SS} \neq T\right) \leq K_1 \exp\left(-K_2 \frac{ac_2 \lambda^2}{L^2}\right). \quad (21)$$

Theorem 5 bounds exponentially the error of the SS algorithm in the case that the contrast function is quite general convex function. It is similar (in fact, slightly worse) than Theorem 3 that considers GLM. Below we discuss its assumptions and compare it in detail with Theorem 3.

**Remark 6** . *The main assumptions of Theorem 5 are convexity of the contrast function and conditions (18) and (19). They can be proved using tools from the empirical process theory such that concentration inequalities (Massart, 2000), the Symmetrization (van der Vaart and Wellner, 1996, Lemma 2.3.1) and Contraction Lemma (Ledoux and Talagrand, 1991, Theorem 4.12). It is quite remarkable that to get (18) or (19) we need only one new condition. Namely, we need that the contrast function is Lipschitz in the following sense: there exists  $L > 0$  such that for all  $x, y, 0 < r \leq \mathring{\beta}_{\min}$  and  $\beta, \tilde{\beta} \in B_1(r)$*

$$|\phi(\beta^T x, y) - \phi(\tilde{\beta}^T x, y)| \leq L|\beta^T x - \tilde{\beta}^T x|. \quad (22)$$

*This fact follows from Massart (2000, Theorem 9) and Bühlmann and van de Geer (2011, Lemma 14.20). We show in the appendix that to get (19) we need similar argumentation and (22) to be satisfied for all  $x, y, r \in (0, \sqrt{\delta_{t-1}})$ ,  $J : T \subset J, r(X_J) = |J| \leq \bar{t}$  and  $\beta, \tilde{\beta} \in B_{2,J}(r)$ . Notice that logistic and absolute contrast functions satisfy (22) with  $L = 2$  and  $L = 1$ , respectively. Moreover, every convex function is locally Lipschitz, so (22) is also satisfied for remaining contrasts (but in these cases  $L$  depends on  $n$ ).*

**Remark 7** *The condition (17) is often called the "margin condition" in the literature. For quadratic and logistic contrasts it has been considered earlier. To prove it for SVM contrasts one can use methods based on the modulus of convexity (Bartlett et al., 2006, Lemma 7).*

**Remark 8 (Comparison to Theorem 3 - Similarities)** *As we have already mentioned Theorem 5 can be also applied to GLM. We can calculate that for quadratic and logistic contrasts we have*

$$\ell(\beta) - \ell(\mathring{\beta}) = -(\beta - \mathring{\beta})^T X^T y + g(\beta) - g(\mathring{\beta})$$

and

$$\mathbb{E}\ell(\beta) - \mathbb{E}\ell(\mathring{\beta}) = -(\beta - \mathring{\beta})^T X^T \mathbb{E}y + g(\beta) - g(\mathring{\beta}),$$

where  $g$  is a total cumulant function. Therefore, the condition (17) for  $l_2$ -balls is the same as (7). Besides, we have for  $\varepsilon = y - \mathbb{E}y$  that

$$\ell(\beta) - \mathbb{E}\ell(\beta) - [\ell(\mathring{\beta}) - \mathbb{E}\ell(\mathring{\beta})] = (\beta - \mathring{\beta})^T X^T \varepsilon.$$

Therefore, we can calculate that for GLM we have  $Z(r)/r \leq |X^T \varepsilon|_\infty$  and  $U_J(r)/r \leq \sqrt{\varepsilon^T H_J \varepsilon}$ , that simplifies bounds in (18) and (19). It makes them analogous to those in Lemma 10 (i) and (iii). Previously they have been obtained using properties of errors  $\varepsilon_1, \dots, \varepsilon_n$ . Now they are proved using the empirical process theory in more general framework. Lemma 10 (i) and (iii) are intended to GLM, so they should give better results than the general methods. And they do, but only in relation to constants  $K_i$ . Finally, notice that the left-hand side in condition (20) and the result in (21) are the same (again with respect to the constants) to those in Theorem 3.

**Remark 9 (Comparison to Theorem 3 - Differences)** *The SS algorithm consists of two steps. In the last paragraph we have mentioned that the theoretical analysis of the second step (selection) is similar in GLM and models with convex contrasts. However, we can find differences while investigating the first step (screening based on the lasso). In general, it is related to the fact that the properties of the lasso in GLM (so we work with differentiable contrasts that can be nicely decomposed) are better and/or better studied. In Theorem 5 we assume (17) with also respect to  $l_1$ -balls, which makes the right-hand side of (20) usually worse than in (12), because  $c_1^2$  appears in (20).*

Summarizing, using the empirical process theory we are able to prove that the SS algorithm works in a satisfactory way even in quite general models with convex contrasts.

## 5. Proofs

In the following subsections we present auxiliary exponential inequalities for subgaussian random variables, upper-bound selection error of TL, and two parts of selection error of SS.

### 5.1. Exponential Bounds for Subgaussian Vectors

In the following lemma we develop auxiliary probabilistic tools. Specifically, in lemma 10 (iii) we generalize Wallace inequality for  $\chi^2$  distribution Wallace (1959) to the subgaussian case using the inequality for the moment generating function in lemma 10 (ii). The last inequality is proved by the decoupling technique as in the proof of Theorem 2.1 in Hsu et al. (2012).

**Lemma 10** *Let  $\varepsilon \in \mathbb{R}^n$  be a vector of zero-mean independent errors having subgaussian distribution with a constant  $\sigma$ ,  $\nu \in \mathbb{R}^n$ ,  $0 < a < 1$  and  $H$  be a orthogonal projection such that  $\text{tr}(H) = m$ . Then*

(i) for  $\tau > 0$

$$\mathbb{P}(\varepsilon^T \nu / \|\nu\| \geq \tau) \leq \exp\left(-\frac{\tau^2}{2\sigma^2}\right) \quad (23)$$

(ii)

$$\mathbb{E} \exp\left(\frac{a}{2\sigma^2} \varepsilon^T H \varepsilon\right) \leq \exp\left(-\frac{m}{2} \log(1-a)\right)$$

(iii) for  $\tau > 1$

$$\mathbb{P}(\varepsilon^T H \varepsilon \geq m\sigma^2\tau) \leq \exp\left(-\frac{m}{2}(\tau - 1 - \log \tau)\right) \quad (24)$$

**Proof** [Proof of Lemma 10] Let  $Z = \varepsilon^T \nu / \|\nu\|$  and  $a > 0$ . From Markov inequality we obtain

$$\mathbb{P}(Z \geq \tau) \leq e^{-a} \mathbb{E} e^{aZ} \leq e^{-a\tau + a^2\sigma^2/2}.$$

Minimizing the last expression w.r.t.  $a$  gives part (i).

Let  $\xi \sim N(0, I_n)$  be a vector of iid standard normal errors independent of  $\varepsilon$ . We have

$$\begin{aligned} \mathbb{E} \exp\left(\frac{a}{2\sigma^2} \varepsilon^T H \varepsilon\right) &= \mathbb{E} \mathbb{E}\left(\exp\left(\frac{\sqrt{a}}{\sigma} \xi^T H \varepsilon\right) \middle| H \varepsilon\right) = \mathbb{E} \exp\left(\frac{\sqrt{a}}{\sigma} \xi^T H \varepsilon\right) \\ &= \mathbb{E} \mathbb{E}\left(\exp\left(\frac{\sqrt{a}}{\sigma} \xi^T H \varepsilon\right) \middle| \xi^T H\right) \leq \mathbb{E} \exp\left(\frac{a}{2} \xi^T H \xi\right). \end{aligned}$$

Thus part (ii) follows from a known formula for the moment generating function of the  $\chi^2$  distribution.

From Markov inequality and part (ii) of the Lemma we have

$$\mathbb{P}(\varepsilon^T H \varepsilon \geq m\sigma^2\tau) \leq \exp\left(-\frac{am\tau}{2}\right) \mathbb{E} \exp\left(\frac{a}{2\sigma^2} \varepsilon^T H \varepsilon\right) \leq \exp\left(-\frac{m}{2}\left(a\tau + \log(1-a)\right)\right).$$

Thus after minimization the last expression w.r.t  $a$  we obtain part (iii).  $\blacksquare$

## 5.2. Proof of Lemma 1

Let  $\mathcal{A}_a = \{|X^T \varepsilon|_\infty \leq a\lambda\}$  and  $\hat{\varepsilon} = y - \dot{\gamma}(X\hat{\beta})$ . We have  $\dot{\ell}(\hat{\beta}) = -X^T \hat{\varepsilon}$  and from the Karush-Kuhn-Tucker (KKT) theorem we obtain equations

$$X^T \hat{\varepsilon} = \lambda [\mathbb{I}(\hat{\beta} > 0) - \mathbb{I}(\hat{\beta} < 0) + u\mathbb{I}(\hat{\beta} = 0)] \quad \text{for } u \in [-1, 1].$$

Let  $\Delta = \hat{\beta} - \check{\beta}$  and  $\nu \in \mathbb{R}^p$  be such that  $\text{sgn}(\nu_{\bar{T}}) = \text{sgn}(\Delta_{\bar{T}})$ . We have  $\nu_J^T X_J^T \hat{\varepsilon} = \lambda |\nu_J|_1$  for  $J \subseteq \bar{T}$  and consequently

$$\begin{aligned} D(\nu) &= \nu^T X^T [\dot{\gamma}(X\hat{\beta}) - \dot{\gamma}(X\check{\beta})] = \nu_T^T X_T^T (\varepsilon - \hat{\varepsilon}) + \nu_{\bar{T}}^T X_{\bar{T}}^T (\varepsilon - \hat{\varepsilon}) \\ &\leq |\nu_T|_1 (|X_T^T \varepsilon|_\infty + |X_{\bar{T}}^T \hat{\varepsilon}|_\infty) + |\nu_{\bar{T}}|_1 (|X_{\bar{T}}^T \varepsilon|_\infty - \lambda) \leq |\nu_T|_1 (1+a)\lambda + |\nu_{\bar{T}}|_1 (a-1)\lambda. \end{aligned} \quad (25)$$

Then letting  $\nu = \Delta_J$  for  $J \subseteq \bar{T}$  we have  $D(\nu) \leq 0$ . Moreover, for  $\nu = \Delta$  we have from convexity of  $g$  that

$$D(\nu) = (\hat{\beta} - \check{\beta})^T [\dot{g}(\hat{\beta}) - \dot{g}(\check{\beta})] \geq 0.$$

Indeed,  $D_0(\beta_1, \beta_2) = (\beta_1 - \beta_2)^T [\dot{g}(\beta_1) - \dot{g}(\beta_2)]$  is the symmetrized Bregman divergence (Huang and Zhang, 2012). Hence  $(1-a)|\nu_{\bar{T}}|_1 \leq (1+a)|\nu_T|_1$ . Thus, on  $\mathcal{A}_a$ ,  $\Delta \in \mathcal{C}_a$  and from the definition of  $\zeta_a$  we obtain using KKT again

$$\zeta_{a,q} |\Delta|_q \leq |X^T [\dot{\gamma}(X\hat{\beta}) - \dot{\gamma}(X\check{\beta})]|_\infty \leq |X^T \hat{\varepsilon}|_\infty + |X^T \varepsilon|_\infty \leq (1+a)\lambda.$$

## 5.3. Proof of Theorem 2

First we will prove that  $\mathcal{A}_a \subseteq \{\hat{T}_{TL} = T\}$ . From Lemma 1 and assumptions we have on  $\mathcal{A}_a$

$$|\Delta|_\infty \leq (1+a)\lambda\zeta_a^{-1} \leq \tau < \check{\beta}_{min}/2. \quad (26)$$

Thus using (26) twice we have for  $j \in T$  and  $k \notin T$

$$|\hat{\beta}_j| \geq |\check{\beta}_j| - |\hat{\beta}_j - \check{\beta}_j| > \check{\beta}_{min} - \check{\beta}_{min}/2 > \tau \geq |\hat{\beta}_k - \check{\beta}_k| = |\hat{\beta}_k| \quad (27)$$

and it follows that  $\mathcal{A}_a \subseteq \{\hat{T}_{TL} = T\}$ . Moreover, the assumptions of the Theorem imply

$$-a_1^2 \lambda^2 + 2\sigma^2 \log p \leq -(1-a_2)a_1^2 \lambda^2.$$

Hence, using Lemma 10 (i) we easily obtain

$$\mathbb{P}(\hat{T}_{TL} \neq T) \leq \mathbb{P}(\mathcal{A}_{a_1}^c) \leq 2p \exp\left(-\frac{a_1^2 \lambda^2}{2\sigma^2}\right) \leq 2 \exp\left(-\frac{(1-a_2)a_1^2 \lambda^2}{2\sigma^2}\right).$$

### 5.4. Proof of Theorem 3

Let us observe that the consecutive steps of SS constitute decomposition of the selection error into 2 parts:  $\{\hat{T} \neq T\} = \{T \notin \mathcal{J}\} \cup \{T \in \mathcal{J}, \hat{T} \neq T\}$ . Therefore Theorem 3 follows easily from (28) and (30) below.

Having in mind that for given  $a_1 \in (1/2, 1)$  we let  $a_2 = 1 - (1 - \log(1 - a_1))(1 - a_1)$ ,  $a_3 = 2 - 1/a_1$  and  $a_4 = \sqrt{a_1 a_2}$ , by arguments similar to those in Theorem 2 we obtain  $\mathcal{A}_{a_4} \subseteq \{T \in \mathcal{J}\}$ . Moreover assumptions  $0 < c \leq 1$  and  $\frac{2\sigma^2 \log p}{a_3 a_2 a_1 c} \leq \lambda^2$  imply

$$-a_4^2 \lambda^2 + 2\sigma^2 \log p \leq -a_2 a_1 c \lambda^2 + 2\sigma^2 \log p \leq -(1 - a_3) a_2 a_1 c \lambda^2 = -a_2 (1 - a_1) c \lambda^2.$$

As a result

$$\mathbb{P}(T \notin \mathcal{J}) \leq \mathbb{P}(\mathcal{A}_{a_4}^c) \leq 2p \exp\left(-\frac{a_4^2 \lambda^2}{2\sigma^2}\right) \leq 2 \exp\left(-\frac{a_2 (1 - a_1) c \lambda^2}{2\sigma^2}\right). \quad (28)$$

Now we bound probability  $\mathbb{P}(T \in \mathcal{J}, \hat{T} \neq T)$ . Components of the selection error set are included in the critical sets of the following form  $\mathcal{C}_J(\tau) = \{\varepsilon^T H_J \varepsilon \geq \tau\}$ .

Proofs of the lemmas stated below are relegated to the Supplemental Materials.

**Lemma 11** *If for  $a \in (0, 1)$  we have  $\lambda^2 < c\delta / (1 + \sqrt{2a})^2$ , then  $\{T \in \mathcal{J}, \hat{T} \subset T\} \subseteq \mathcal{C}_T(ac\lambda^2)$ .*

**Lemma 12** *For  $a \in (0, 1)$  we have*

$$\{T \in \mathcal{J}, \hat{T} \supset T\} \subseteq \mathcal{C}_T((1 - a)c\lambda^2) \cup \bigcup_{J \supset T} \mathcal{C}_{J \setminus T}(|J \setminus T|ac\lambda^2).$$

Let us define  $\tau_0 = \frac{1}{1 - a_1}$ ,  $\tau_1 = \frac{(1 - a_1)c\lambda^2}{t\sigma^2}$  and  $\tau_2 = \frac{a_1 c \lambda^2}{\sigma^2}$ . Under our assumptions we have  $2 < \tau_0 < \tau_1 < \tau_2$ . Let  $f_2(\tau) = 1 - (1 + \log \tau)/\tau$  for  $\tau > 1$ . Of course  $f_2$  is increasing,  $f_2(1) = 0$  and  $f_2(\tau) \rightarrow 1$  for  $\tau \rightarrow \infty$ . Consequently  $a_2 = f_2(\tau_0) < f_2(\tau_1) < f_2(\tau_2)$ , which means that

$$a_2 \tau_r < \tau_r - 1 - \log \tau_r \text{ for } r = 1, 2. \quad (29)$$

From Lemma 11, Lemma 12, Lemma 10 (iii) and (29) we get

$$\begin{aligned} \mathbb{P}(T \in \mathcal{J}, \hat{T} \neq T) &\leq \mathbb{P}(\mathcal{C}_T(t\sigma^2 \tau_1)) + \sum_{J \supset T} \mathbb{P}(\mathcal{C}_{J \setminus T}(|J \setminus T| \sigma^2 \tau_2)) \leq \exp(-ta_2 \tau_1 / 2) \\ &+ \sum_{m=1}^{p-t} \binom{p-t}{m} \exp(-ma_2 \tau_2 / 2) \leq \exp(-ta_2 \tau_1 / 2) + \sum_{m=1}^{p-t} \frac{1}{m!} \exp\left(\frac{-m}{2}(a_2 \tau_2 - 2 \log p)\right). \end{aligned}$$

Using  $\exp(d) - 1 \leq \log(2)^{-1} d$  for  $0 \leq d \leq \log(2)$  and the fact that probability is not greater than 1 we obtain

$$\mathbb{P}(T \in \mathcal{J}, \hat{T} \neq T) \leq \exp(-ta_2 \tau_1 / 2) + (\log 2)^{-1} \exp(-(a_2 \tau_2 - 2 \log p) / 2).$$

For  $a_3 \in (0, 1)$ , assumption  $\frac{2\sigma^2 \log p}{a_3 a_2 a_1 c} \leq \lambda^2$  implies  $-a_2 \tau_2 + 2 \log p \leq -(1 - a_3) a_2 \tau_2$ , therefore

$$\mathbb{P}(T \in \mathcal{J}, \hat{T} \neq T) \leq \exp(-ta_2 \tau_1 / 2) + (\log 2)^{-1} \exp(-(1 - a_3) a_2 \tau_2 / 2). \quad (30)$$

## 6. Experiments

While convenient for theoretical analysis TL, FCP, SOS or SS algorithms use the Lasso estimators only for one value of the penalty, the practical Lasso implementations return coefficient estimators for a given net of it (R package `glmnet` described in [Friedman et al. \(2010\)](#)). Similarly, using a net of penalty values, the Minimax Concave Penalty (MCP) algorithm, a popular realization of FCP, has been implemented for linear (in R package `SparseNet` [Mazumder et al. \(2011\)](#)) and logistic models (in R package `cvplogistic` [Jiang and Huang \(2014\)](#)).

In order to improve SOS and SS performance, we propose a net modification of SOS called SOSnet (Algorithm 2 below), which employs estimates for  $m$  values:  $\lambda_1, \dots, \lambda_m$ . This alteration results in higher accuracy of model selection and sparse prediction as shown in our experiments. The most intensive step of the algorithm, namely computations of Lasso estimators, can be performed for all  $\lambda$  values during one run of the `glmnet` algorithm. As in SOS the ordering step allows to choose  $T$  when the screening step is correct i.e.  $\text{supp}\hat{\beta}_\lambda \supseteq T$  but ordering given by absolute values of coordinates of  $\hat{\beta}_\lambda$  is wrong. An additional loop (for  $l = 1$  series to  $o$ ) is introduced in order to find a correct screening set having possibly small cardinality.

---

### Algorithm 2 SOSnet

---

**Input:**  $y, X$  and  $(o, \lambda \leq \lambda_1 < \dots < \lambda_m)$   
**Screening** (Lasso)  
**for**  $k = 1$  **to**  $m$  **do**  
 $\hat{\beta}^{(k)} = \text{argmin}_\beta \{\ell(\beta) + \lambda_k |\beta|_1\}$ ; order nonzero  $|\hat{\beta}_{j_1}^{(k)}| \geq \dots \geq |\hat{\beta}_{j_{s_k}}^{(k)}|$ , where  $s_k = |\text{supp}\hat{\beta}^{(k)}|$ ;  
**Ordering** (squared Wald tests)  
**for**  $l = 1$  **to**  $o$  **do**  
set  $J = \{j_1, j_2, \dots, j_{s_{kl}}\}$ , where  $s_{kl} = \lfloor \frac{s_k \cdot l}{o} \rfloor$ ;  
compute  $\hat{\beta}_J^{ML}$ ;  
set predictors in  $J$  according to squared Wald tests:  $w_{i_1}^2 \geq w_{i_2}^2 \geq \dots \geq w_{i_{s_{kl}}}^2$  ;  
set  $\mathcal{J}_{kl} = \{\{i_1\}, \{i_1, i_2\}, \dots, \{i_1, i_2, \dots, i_{s_{kl}}\}\}$   
**end for**;  
**end for**;  
**Selection** (GIC)  
 $\mathcal{J} = \bigcup_{k=1}^m \bigcup_{l=1}^o \mathcal{J}_{kl}$ ;  $\hat{T} = \text{argmin}_{J \in \mathcal{J}} \{\ell_J + \lambda^2/2|J|\}$   
**Output:**  $\hat{T}, \hat{\beta}^{GSSG} = \hat{\beta}_{\hat{T}}^{ML}$

---

We performed numerical experiments fitting sparse linear and logistic models to high-dimensional benchmark simulations and real data sets.

### 6.1. Simulated Data

For linear models we studied the performance of two algorithms: SOSnet and MCP computed using the R package `SparseNet` [Mazumder et al. \(2011\)](#) for the default 9 values of  $\gamma$  and 50 values of  $\lambda$ . Our algorithm used the R package `glmnet` [Friedman et al. \(2010\)](#) to compute the Lasso estimators for 50 lambdas on a log scale and with  $o = 5$ .

Table 1: Plan of experiments for linear models.

	n	p	$\beta$	$\rho$	$\sigma^2$	SNR
N.1.5	100	3000	$\beta^{(1)}$	.5	4	2.3
N.1.7	100	3000	$\beta^{(1)}$	.7	4	2.6
N.1.9	100	3000	$\beta^{(1)}$	.9	4	3
N.2.5	200	2000	$\beta^{(2)}$	.5	7	2.4
N.2.7	200	2000	$\beta^{(2)}$	.7	7	2.3
N.2.9	200	2000	$\beta^{(2)}$	.9	7	2.2

We generated samples  $(y_i, x_i)$ ,  $i = 1, \dots, n$  from the normal linear model. Two vectors of parameters were considered:  $\beta^{(1)} = (3, 1.5, 0, 0, 2, 0_{p-5}^T)^T$ , as in Wang et al. (2013) as well as  $\beta^{(2)} = (0_{p-10}^T, s_1 \cdot 2, s_2 \cdot 2, \dots, s_{10} \cdot 2)^T$ , where  $s_l$  equals 1 or -1 with equal probability,  $l = 1, \dots, 10$  chosen separately for every run as in experiment 2 in Wang et al. (2014). The rows of  $X$  were iid  $p$ -dimensional vectors  $x_i \sim N(0_p, \Xi)$ . We considered auto-regressive structure of covariance matrix that is  $\Xi = (\rho^{|i-j|})_{i,j=1}^p$  for  $\rho = 0.5, 0.7, 0.9$ . The columns of  $X$  were centred and normalized so that  $\|x_{\cdot j}\|^2 = n$  and  $\varepsilon \sim N(0_n, \sigma^2 I_n)$ . The plan of experiments is presented in Table 1 with SNR meaning a *Signal to Noise Ratio*.

For every experiment the results were based on  $N = 1000$  simulation runs. We reported mean model dimension (MD) that is  $|\text{supp}(\hat{\beta})|$  and mean squared prediction error (PE) on new data set with 1000 observations equalling  $\|X\hat{\beta} - X\beta\|^2 / (n\sigma^2)$ . We chose the model using GIC with  $\lambda^2 = c \cdot \log(p) \cdot \sigma^2$ . For each value of hyperparameter  $c = .25, .5, \dots, 7.5$  values of  $(MD(c), PE(c))$  for the models chosen by GIC=GIC( $c$ ) were calculated and averaged over simulations. The results are presented in two first columns of Figure 1. The two vertical lines indicate models chosen using GIC with  $c = 2.5$ : the black one for SOSnet and the red one for SparseNet. The blue vertical line denotes the true model dimension.

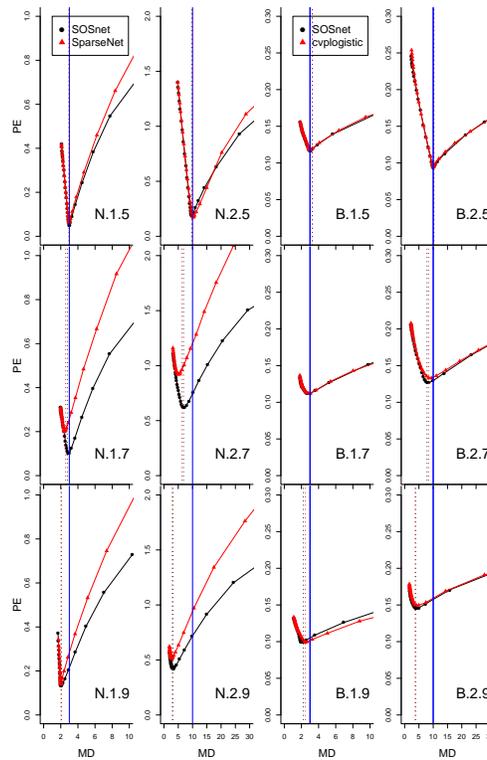
For logistic models we compared the performance of two algorithms: SOSnet and MCP implemented in the R package `cvplogistic` for the default value of  $\kappa = 1/2.7$  and 100 values of  $\lambda$ . As for linear models, SOSnet called the R package `glmnet` Friedman et al. (2010) to compute the Lasso estimators for 20 lambdas on a default log scale and with  $o = 5$ . We performed experiments very similar to those for linear models, changing only  $n$  and the number of simulation runs to  $N = 500$ . The plan of experiments is shown in the Table 2. Random samples were generated according to the binomial distribution. We reported prediction error defined as misclassification frequency on new data set with 1000 observations. The results organized in a similar way as for the linear models are shown in columns 3–4 of Figure 1. The two vertical lines indicate models chosen using GIC with  $c = 2$ , the black one for SOSnet and the red one for `cvplogistic`.

Summarizing the results of the simulation study, one can observe that SOSnet for linear models turned out to have equal or lower PE in almost all of the experimental setups. The differences are most visible in setups with autocorrelation structure with  $\rho = 0.7$ . The value  $c = 2.5$  in GIC usually gave satisfactory results. The mean execution time of SOSnet was approximately 1.5 times longer than for SparseNet. SOSnet for logistic regression gave similar accuracy as `cvplogistic` with much lower execution time: SOSnet was approximately 10 times faster. The value  $c = 2$  in GIC usually gave satisfactory results.

Table 2: Plan of experiments for logistic models.

	n	p	$\beta$	$\rho$
B.1.5	300	3000	$\hat{\beta}^{(1)}$	.5
B.1.7	300	3000	$\hat{\beta}^{(1)}$	.7
B.1.9	300	3000	$\hat{\beta}^{(1)}$	.9
B.2.5	500	2000	$\hat{\beta}^{(2)}$	.5
B.2.7	500	2000	$\hat{\beta}^{(2)}$	.7
B.2.9	500	2000	$\hat{\beta}^{(2)}$	.9

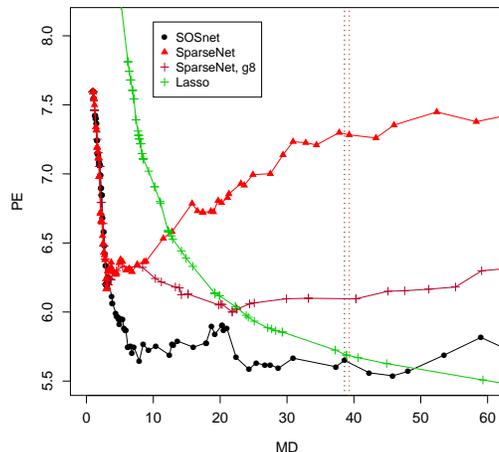
Figure 1: Results for simulated data



## 6.2. Real Data Sets

The methylation data set was described in [Hannum et al. \(2013\)](#). It consist of the age of 656 human individuals together with values of phenotypic features such as gender and body mass index and of genetic features, which are methylation states of 485 577 CpG markers. Methylation was recorded as a fraction representing the frequency of methylation of a given CpG marker across the population of blood cells taken from a single individual. In our comparison we used only genetic features from which we extracted 193 870 most relevant CpGs according to onefold t-tests with Benjamini-Hochberg adjustment, FDR=.05. We compared the root mean squared errors (PE) and model dimensions (MD) for SOSnet and SparseNet via 10-fold cross-validation. For each value of hyperparameter  $c = .25, .5, \dots, 7.5$  values of  $(MD(c), PE(c))$  for the models chosen

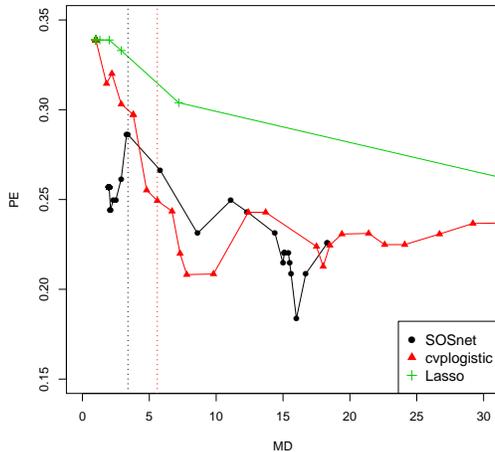
Figure 2: Results for the methylation data set



by  $GIC=GIC(c)$  were calculated and averaged over 10 folds. The results are presented in Figure 2. SparseNet yields a path of models for each value of parameter  $\gamma = g_1, \dots, g_9$ . We present results for  $g_1$ , corresponding to the Lasso, for  $g_9$ , close to the best subset, and for an intermediate value  $g_8$  in Figure 2. Remarkably, SOSnet gives uniformly smaller PE than SparseNet for all  $MD \geq 3$ . The two vertical lines indicate models chosen using GIC with  $c = 2.5$ : the black one for SOSnet and the red one for SparseNet.

A logistic model was fitted to the breast cancer data described in Gravier et al. (2010) which concerns small, invasive carcinomas without axillary lymph node involvement to predict metastasis of small node-negative breast carcinoma. There were 168 patients: 111 with no event after diagnosis labeled as good, and 57 with early metastasis labeled as poor. The number of predictors in this data was 2905. We compared the mean errors of binary prediction (PE) and model dimensions (MD) for SOSnet and `cvplogistic` via 10-fold cross-validation. The results are presented in Figure 3. Minimal PE for SOSnet was smaller than for `cvplogistic`, but for a larger model. The algorithms work comparably, but again SOSnet was 63 times faster. The two vertical lines indicate models chosen using GIC with  $c = 2$ : the black one for SOSnet and the red one for `cvplogistic`.

Figure 3: Results for the breast cancer data set



## 7. Conclusions

In the paper we propose the SS algorithm which is an alternative method to TL and FCP of improving the Lasso. For linear models it seems to be the benchmark for the theory of model selection as it is constructive, computationally efficient and leads to consistent model selection under weak assumptions.

Our approach encompasses fundamental model for prediction of continuous as well as of binary response and the main result is stated jointly for both of them. Its assumptions are stated in the most general form which allows proving exponential bound without obscuring the essence of the results and comparing the bounds for both models. By simplifying SOS to SS we were able to simplify reasoning used for SOS and then extend them from linear models to general predictive models.

We propose an algorithm SOSnet, which is a generalization of the SOS algorithm for general predictive models. Using net of parameters, SOSnet avoids problem of choosing one specific  $\lambda$ . The gap between theoretical results for SS and the SOSnet algorithm is similar to the difference between theory for FCP and its implementations `SparseNet` or `cvplogistic`. Numerical experiments reveal that for linear models SOSnet is more accurate than `SparseNet` with comparable computing time, whereas for logistic models performance of SOSnet is on par with performance of `cvplogistic` with computing times at least 10 times shorter.

We have shown in simulations (dotted vertical lines in Figure 1) that predictively optimal  $\lambda$  for linear models equals approximately  $\sqrt{2.5\sigma^2 \log p}$ , which is close to (14) and for logistic models is  $\sqrt{2 \log p}$ , which together with (15) suggests that  $c \approx 1/4$ . The relations between the safest choice  $\lambda$  discussed in Remark 2 and predictively optimal  $\lambda$  are important applications of our theory.

## Acknowledgments

We thank a bunch of people.

## References

- P. L. Bartlett, M. I. Jordan, and J. D. McAuliffe. Convexity, classification and risk bounds. *Journal of the American Statistical Association*, 101:138–156, 2006.
- P. Bühlmann and S. van de Geer. *Statistics for High-dimensional Data*. Springer, New York, 2011.
- B. Efron, T. Hastie, I. Johnstone, and R. Tibshirani. Least angle regression. *Annals of Statistics*, 32:407–499, 2004.
- J. Fan and R. Li. Variable selection via nonconcave penalized likelihood and its oracle properties. *Journal of the American Statistical Association*, 96:1348–1360, 2001.
- J. Fan, L. Xue, and H. Zou. Strong oracle optimality of folded concave penalized estimation. *Annals of Statistics*, 42:819–849, 2014.
- J. Friedman, T. Hastie, and R. Tibshirani. Regularization paths for generalized linear models via coordinate descent. *Journal of Statistical Software*, 33:1–22, 2010.
- E. Gravier, G. Pierron, A. Vincent-Salomon, N. Gruel, V. Raynal, A. Savignoni, Y. De Rycke, J.Y. Pierga, C. Lucchesi, F. Reyal, et al. Prognostic DNA signature for T1T2 node-negative breast cancer patients. *Genes, Chromosomes and Cancer*, 49(12):1125–1125, 2010.
- G. Hannum, J. Guinney, L. Zhao, L. Zhang, G. Hughes, S. Sada, B. Klotzle, M. Bibikova, J.B. Fan, Y. Gao, et al. Genome-wide methylation profiles reveal quantitative views of human aging rates. *Molecular cell*, 49(2):359–367, 2013.
- D. Hsu, S.M. Kakade, and T. Zhang. A tail inequality for quadratic forms of subgaussian random vectors. *Electronic Communications in Probability*, 17:1–6, 2012.
- J. Huang and C.H. Zhang. Estimation and selection via absolute penalized convex minimization and its multistage adaptive applications. *Journal of Machine Learning Research*, 13:1839–1864, 2012.
- D. Jiang and J. Huang. Majorization minimization by coordinate descent for concave penalized generalized linear models. *Statistics and Computing*, 24(5):871–883, 2014.
- M. Ledoux and M. Talagrand. *Probability in Banach Spaces: Isoperimetry and Processes*. Springer, Berlin, 1991.
- P. Massart. About the constants in talagrand's concentration inequalities for empirical processes. *The Annals of Probability*, 28:863–884, 2000.
- R. Mazumder, J.H. Friedman, and T. Hastie. SparseNet: Coordinate descent with nonconvex penalties. *Journal of the American Statistical Association*, 106:1125–1138, 2011.
- N. Meinshausen and P. Bühlmann. High dimensional graphs and variable selection with the Lasso. *Annals of Statistics*, 34:1436–1462, 2006.
- P. Pokarowski and J. Mielniczuk. Combined  $\ell_0$  and  $\ell_1$  penalized least squares for linear model selection. *Journal of Machine Learning Research*, 16:961–992, 2015.

- X. Shen, W. Pan, and Y. Zhu. Likelihood-based selection and sharp parameter estimation. *Journal of the American Statistical Association*, 107:223–232, 2012.
- X. Shen, W. Pan, Y. Zhu, and H. Zhou. On constrained and regularized high-dimensional regression. *Annals of the Institute of Statistical Mathematics*, 65:807–832, 2013.
- R. Tibshirani. Regression shrinkage and selection via the Lasso: a retrospective. *Journal of the Royal Statistical Society Series B*, 73:273–282, 2011.
- A. W. van der Vaart and J. A. Wellner. *Weak Convergence and Empirical Processes: With Applications to Statistics*. Springer Verlag, New York, 1996.
- D.L. Wallace. Bounds on normal approximations to student’s and the chi-square distributions. *Annals of Mathematical Statistics*, 30:1121–1130, 1959.
- L. Wang, Y. Kim, and R. Li. Calibrating non-convex penalized regression in ultra-high dimension. *Annals of Statistics*, 41:2505–2536, 2013.
- Z. Wang, H. Liu, and T. Zhang. Optimal computational and statistical rates of convergence for sparse nonconvex learning problems. *Annals of Statistics*, 42:2164–2201, 2014.
- F. Ye and C.H. Zhang. Rate minimaxity of the Lasso and Dantzig Selector for the lq loss in lq balls. *Journal of Machine Learning Research*, 11:3519–3540, 2010.
- C.H. Zhang. Nearly unbiased variable selection under minimax concave penalty. *Annals of Statistics*, 38:894–942, 2010a.
- C.H. Zhang and T. Zhang. A general theory of concave regularization for high-dimensional sparse estimation problems. *Statistical Science*, 27:576–593, 2012.
- T. Zhang. Analysis of multi-stage convex relaxation for sparse regularization. *Journal of Machine Learning Research*, 11:1081–1107, 2010b.
- P. Zhao and B. Yu. On model selection consistency of Lasso. *Journal of Machine Learning Research*, 7:2541–2563, 2006.
- S. Zhou. Thresholding procedures for high dimensional variable selection and statistical estimation. In *NIPS*, pages 2304–2312, 2009.

## Appendix A. Proofs and auxiliary results

Let us define for  $k = 1, \dots, t - 1$

$$\begin{aligned} \mathcal{E}_k(\tau) &= \{\exists J \subset T \mid |T \setminus J| = k : \ell_J - \ell_T \leq \tau\}, \\ \delta_k &= \min_{J \subset T, |T \setminus J| = k} \|(I - H_J)X_T \hat{\beta}_T\|^2, \\ B_k &= \{\beta_T : \|X_T \beta_T - X_T \hat{\beta}_T\|^2 \leq \delta_k\}. \end{aligned}$$

Obviously  $\delta = \min_k \delta_k/k$ . First we prove the following lemma which will be used in the proof of Lemma 3.

**Lemma 5.** For  $b \in (0, 1)$  we have

$$\mathcal{E}_k(bc\delta_k/2) \subseteq \mathcal{C}_T((1-b)^2c^2\delta_k/4).$$

**Proof** [Proof of Lemma 5] For  $\beta_T \in \partial B_k$  from assumption (5), Schwartz inequality and properties of the orthogonal projection  $H_T$  we get

$$\begin{aligned} \dot{\ell}(\beta_T) := \ell(\beta_T) - \ell(\hat{\beta}_T) &\geq (\hat{\beta}_T - \beta_T)^T X_T^T H_T \varepsilon + \frac{c}{2} (\hat{\beta}_T - \beta_T)^T X_T^T X_T (\hat{\beta}_T - \beta_T) \\ &\geq -\sqrt{\delta_k \varepsilon^T H_T \varepsilon} + \frac{c}{2} \delta_k. \end{aligned}$$

Since the last expression does not depend on  $\beta_T$ , we have for  $b \in (0, 1)$

$$\mathcal{L}_k(b) = \left\{ \min_{\beta_T \in \partial B_k} \dot{\ell}(\beta_T) \leq \frac{bc\delta_k}{2} \right\} \subseteq \left\{ -\sqrt{\delta_k \varepsilon^T H_T \varepsilon} + \frac{c}{2} \delta_k \leq \frac{bc\delta_k}{2} \right\} = \mathcal{C}_T \left( \frac{(1-b)^2c^2\delta_k}{4} \right).$$

Let us notice that for  $J \subset T$  such that  $|T \setminus J| = k$  we have  $\|X_T(\hat{\beta}_J^{ML} - \hat{\beta}_T)\|^2 \geq \delta_k$ , so  $\hat{\beta}_J^{ML} \notin \text{int}(B_k)$ . Since  $\dot{\ell}$  is convex and  $\dot{\ell}(\hat{\beta}_T) = 0$  we obtain  $\mathcal{L}_k(b) \supseteq \mathcal{E}_k(bc\delta_k/2)$ . ■

**Proof** [Proof of Lemma 3] From assumption of the Lemma  $\lambda^2 < c\delta$ , so  $b_k = \lambda^2 k / (c\delta_k) < 1$ . Hence for  $k = 1, \dots, t-1$

$$k\lambda^2 \leq b_k c \delta_k. \quad (31)$$

Moreover, if  $a \in (0, 1)$  then

$$akc\lambda^2 \leq (1-b_k)^2 c^2 \delta_k / 4, \quad (32)$$

because  $ab \leq (1-b)^2/4$  for  $b = \lambda^2/(c\delta)$ . The last inequality is true if

$$b \leq 1 + 2a - \sqrt{(1+2a)^2 - 1} = f_1(a). \quad (33)$$

Indeed, it is easy to check that (33) follows from the assumption as

$$f_1(a) = \frac{1}{1+2a + \sqrt{(1+2a)^2 - 1}} \geq \frac{1}{(1+\sqrt{2a})^2}.$$

Finally from (31), Lemma 5 and (32) we obtain, respectively

$$\{T \in \mathcal{J}, \hat{T} \subset T\} \subseteq \bigcup_{k=1}^{t-1} \mathcal{E}_k\left(\frac{k\lambda^2}{2}\right) \subseteq \bigcup_{k=1}^{t-1} \mathcal{E}_k\left(\frac{b_k c \delta_k}{2}\right) \subseteq \bigcup_{k=1}^{t-1} \mathcal{C}_T\left(\frac{(1-b_k)^2 c^2 \delta_k}{4}\right) \subseteq \mathcal{C}_T(ac\lambda^2)$$

**Proof** [Proof of Lemma 4] For  $J \supset T$  define  $W_J = X_J(\hat{\beta}_J - \hat{\beta}_J^{ML})$  and  $m = |J \setminus T|$ . We have from assumption (5) and properties of orthogonal projection  $H_J$

$$\begin{aligned} \ell_J - \ell_T &\geq \ell(\hat{\beta}_J^{ML}) - \ell(\hat{\beta}_J) \geq (\hat{\beta}_J - \hat{\beta}_J^{ML})^T X_J^T H_J \varepsilon + \frac{c}{2} (\hat{\beta}_J - \hat{\beta}_J^{ML})^T X_J^T X_J (\hat{\beta}_J - \hat{\beta}_J^{ML}) \\ &= W_J^T H_J \varepsilon + \frac{c}{2} W_J^T W_J = \frac{1}{2c} \|cW_J + H_J \varepsilon\|^2 - \frac{1}{2c} \varepsilon^T H_J \varepsilon \geq -\frac{1}{2c} \varepsilon^T H_J \varepsilon, \end{aligned}$$

so  $\{\ell_J - \ell_T \leq -m\lambda^2/2\} \subseteq \{\varepsilon^T H_J \varepsilon \geq mc\lambda^2\}$ .

Moreover,  $\varepsilon^T H_J \varepsilon \leq \varepsilon^T H_T \varepsilon + \varepsilon^T H_{J \setminus T} \varepsilon$ , hence we obtain for  $\tau > 0$  and  $a \in (0, 1)$

$$\{\varepsilon^T H_J \varepsilon \geq \tau\} \subseteq \{\varepsilon^T H_T \varepsilon \geq (1-a)\tau\} \cup \{\varepsilon^T H_{J \setminus T} \varepsilon \geq a\tau\}.$$

Finally

$$\{T \in \mathcal{J}, \hat{T} \supset T\} \subseteq \bigcup_{J \supset T} \{\ell_J - \ell_T \leq -|J \setminus T|\lambda^2/2\} \subseteq \mathcal{C}_T((1-a)c\lambda^2) \cup \bigcup_{J \supset T} \mathcal{C}_{J \setminus T}(|J \setminus T|ac\lambda^2).$$

■