

Non-Asymptotic Confidence Regions for Regularized Linear Regression Estimates

Balázs Csanád Csáji

Abstract Building *confidence regions* for regression models is of high importance, for example, they can be used for uncertainty quantification and are also fundamental for robust optimization. In practice, these regions are often computed from the asymptotic distributions, which however only lead to heuristic confidence sets. Sign-Perturbed Sums (SPS) is a resampling method which can construct *exact, non-asymptotic, distribution-free* confidence regions under very mild statistical assumptions. In its standard form, the SPS regions are built around the least-squares estimate of linear regression problems, and have favorable properties, such as they are star convex, strongly consistent, and have efficient ellipsoidal outer-approximations. In this paper, we extend the SPS method to *regularized* estimates, particularly, we present variants of SPS for ridge regression, LASSO and elastic net regularization.

1 Introduction

Estimating models based on noisy measurements is a fundamental problem for many scientific, engineering and economic applications. A very important issue in practice is to quantify the uncertainty of the obtained models. This is often done by building confidence regions for the models. While these regions are frequently built using the limiting distribution of the used point-estimate [6], such regions are not guaranteed for finite samples, and can only be seen as heuristics. It is of high importance to construct confidence regions with non-asymptotic guarantees, using minimal statistical assumptions. Resampling methods, such as bootstrap and Monte Carlo approaches, typically use some regularity of the noise to build such regions.

Sign-Perturbed Sums (SPS) is a recently developed resampling method with favorable properties. SPS can construct *exact*, distribution-free confidence regions for

Balázs Csanád Csáji

EPIC Centre of Excellence, MTA SZTAKI: Institute for Computer Science and Control, Hungarian Academy of Sciences, Budapest, Hungary, e-mail: balazs.csaji@sztaki.mta.hu

finite samples [2, 7]. Its standard form constructs (star convex, strongly consistent) confidence sets around the least-squares estimate of linear regression problems.

Regularization is an important tool in regression which can help, for example, to handle ill-posed and ill-conditioned problems, reduce over-fitting, enforce sparsity, and in general to control the shape and smoothness of the regression function. The paper extends SPS to various regularized linear regression problems, particularly, to ridge regression (Tikhonov regularization), LASSO and elastic net regularization.

2 Preliminaries: Asymptotic Confidence Ellipsoids

We start by recalling the standard “textbook” approach to build (asymptotic) confidence ellipsoids around the least-squares estimate of linear regression problems.

Assume we are given a data sample, $\mathcal{D}_n \doteq \{(\varphi_1, y_1), \dots, (\varphi_n, y_n)\}$, with

$$y_t \doteq \varphi_t^T \theta^* + \varepsilon_t, \quad \text{for } t = 1, \dots, n \quad (1)$$

where y_t is the *output*, φ_t is the *input* or *regressor* and ε_t is the (non-observable) *noise* for measurement t . We aim at estimating the (constant) “true” parameter, $\theta^* \in \mathbb{R}^d$. We assume that $\{\varphi_t\} \subset \mathbb{R}^d$ are *deterministic* and the noise $\{\varepsilon_t\}$ is an *independent* sequence of random variables, each having a *symmetric* distribution about zero, that is the distribution of ε_t is the same as that of $-\varepsilon_t$. Finally, for simplicity, we assume that the *regressor matrix*, $\Phi \doteq [\varphi_1, \dots, \varphi_n]^T$, is skinny ($n > d$) and full rank.

One of the standard estimators is the well-known *least-squares* (LS) method

$$\hat{\theta}_n \doteq \arg \min_{\theta \in \mathbb{R}^d} V(\theta | \mathcal{D}_n) = \arg \min_{\theta \in \mathbb{R}^d} \frac{1}{2} \|y - \Phi \theta\|_2^2, \quad (2)$$

where $y \doteq [y_1, \dots, y_n]^T$; $\hat{\theta}_n$, can be obtained from the *normal equation*, that is

$$\nabla_{\theta} V(\hat{\theta}_n | \mathcal{D}_n) = \Phi^T \Phi \hat{\theta}_n - \Phi^T y = 0, \quad (3)$$

which has a unique analytical solution, famously given by $\hat{\theta}_n = (\Phi^T \Phi)^{-1} (\Phi^T y)$.

A crucial question is that how can we *quantify the uncertainty* of the so obtained estimate? This question can be answered, e.g., by constructing *confidence regions* around the point-estimate. More precisely, given a confidence probability $p \in (0, 1)$, we aim at finding a region, $\hat{\Theta}_{\mathcal{D}_n, p}$ around $\hat{\theta}_n$, such that $\mathbb{P}(\theta^* \in \hat{\Theta}_{\mathcal{D}_n, p}) \geq p$.

The standard method to build such regions is to use the *asymptotic distribution* of the estimate [6]. It is known that the (scaled) error of LS is asymptotically Gaussian,

$$\sqrt{n}(\hat{\theta}_n - \theta^*) \xrightarrow{d} \mathcal{N}(0, \sigma^2 R^{-1}), \quad \text{as } n \rightarrow \infty, \quad (4)$$

where $\mathcal{N}(\mu, \Sigma)$ is the (multivariate) Gaussian distribution with mean μ and covariance Σ . This property holds under various conditions, e.g., if the regressors are bounded, there exists a positive definite matrix R as the limit of matrices $R_n \doteq \frac{1}{n} \Phi_n^T \Phi_n$, and $\{\varepsilon_t\}$ are i.i.d. as well as $\mathbb{E}[\varepsilon_t] = 0$ and $\mathbb{E}[\varepsilon_t^2] = \sigma^2$, with $0 < \sigma^2 < \infty$.

Using the limiting distribution, a (heuristic) *confidence ellipsoid* can be built by

$$\tilde{\Theta}_{n,p} \doteq \left\{ \theta \in \mathbb{R}^d : (\theta - \hat{\theta}_n)^T R_n (\theta - \hat{\theta}_n) \leq \frac{q \hat{\sigma}_n^2}{n} \right\}, \quad (5)$$

where $p = F_{\chi^2(d)}(q)$, with $F_{\chi^2(d)}$ being the CDF of the χ^2 distribution with d degrees of freedom; and $\hat{\sigma}_n^2$ is an (unbiased) estimate of the noise variance, that is

$$\hat{\sigma}_n^2 \doteq \frac{1}{n-d} \sum_{t=1}^n (y_t - \varphi_t^T \hat{\theta}_n)^2. \quad (6)$$

Then, we *approximately* have $\mathbb{P}(\theta^* \in \tilde{\Theta}_{n,p}) \approx p$ (and, obviously, $\hat{\theta}_n \in \tilde{\Theta}_{n,p}$).

However, the confidence regions constructed using the asymptotic distribution are *not guaranteed* for finite samples, and are typically *imprecise* if the sample size is small. Another drawback of the asymptotic approach is that it presupposes the *existence* of a limiting distribution, which cannot be guaranteed in certain cases.

3 Sign-Perturbed Sums: Non-Asymptotic Confidence Regions

Now, we overview the *Sign-Perturbed Sums* (SPS) method [2, 7] that can construct *exact, non-asymptotic, distribution-free* confidence regions around the LS estimate.

As first glance, SPS can be seen as a *hypothesis testing* method. It tests the null hypothesis $\theta = \theta^*$, against the alternative hypothesis $\theta \neq \theta^*$. SPS is based on the idea that if $\theta = \theta^*$, then (1) we can compute the exact realization of the noise vector, $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)^T$, by “inverting” the system, and (2) using some *regularity* of the noise (e.g., symmetry), alternative noise realizations can be generated, leading to alternative samples and estimates, which behave “similarly” (in the statistical sense) to the original ones. On the other hand, if $\theta \neq \theta^*$, then the residuals will be biased and the alternative samples and estimates based on them will behave statistically differently than the original ones. SPS applies a rank-test to decide whether the perturbed objects are similar to the original ones. Unlike other resampling based approaches, SPS avoids actually constructing the alternative samples and fitting surrogate models to them, as it directly perturbs the *gradient* of the objective function.

The principal building blocks of SPS are the following *evaluation functions*,

$$Z_i(\theta) \doteq \|\Psi^{1/2} \Phi^T G_i (y - \Phi \theta)\|_2^2, \quad (7)$$

for $i \in \{0, 1, \dots, m-1\}$, where $\Psi = (\Phi^T \Phi)^{-1}$, $m > 0$ is a user-chosen integer, $G_0 \doteq I$, the identity matrix, and for $i \neq 0$, $G_i \doteq \text{diag}(\alpha_{i,1}, \dots, \alpha_{i,n})$; $\{\alpha_{i,j}\}$ are i.i.d. Rademacher variables¹; and $\text{diag}(\cdot)$ builds a diagonal matrix from its argument.

Notice that, apart from an (optional) linear transformation, $\Psi^{1/2}$, whose role is to make a covariance correction, $Z_0(\theta)$ is basically the norm of the (negative) *gradient* of the least-squares objective. The difference between $Z_0(\theta)$ and $Z_i(\theta)$, $i \neq 0$, is that in latter functions the signs of the residuals $(y - \Phi \theta)$ are perturbed in the gradient.

¹ Random variables which take values $+1$ and -1 with probability $1/2$ each.

In case $\theta = \theta^*$, the residuals are the true noises, $y - \Phi\theta^* = \varepsilon$, and we know from the symmetry assumption that for all i , ε and $G_i\varepsilon$ have the same distribution, where G_i is a diagonal matrix containing random signs as defined above. Then,

$$Z_0(\theta^*) = \|\Psi^{1/2}\Phi^T\varepsilon\|_2^2 \stackrel{d}{=} \|\Psi^{1/2}\Phi^TG_i\varepsilon\|_2^2 = Z_i(\theta^*), \quad (8)$$

for $i = 1, \dots, m-1$, where “ $\stackrel{d}{=}$ ” denotes equality in distribution. Nevertheless, variables $\{Z_i(\theta^*)\}$ are of course not independent. On the other hand, it can be proved [2] that they are *conditionally i.i.d.*, conditioned on the σ -algebra generated by $\{|\varepsilon_i|\}$. Consequently, they are also *exchangeable* and hence each ordering² of them, $Z_{i_0}(\theta^*) \prec \dots \prec Z_{i_{m-1}}(\theta^*)$, has the same probability, namely, $1/m!$.

If however, $\theta \neq \theta^*$, then this exchangeability argument does not hold, moreover, $Z_0(\theta)$ will eventually dominate $\{Z_i(\theta)\}_{i \neq 0}$ with high probability as $\|\theta - \theta^*\| \rightarrow \infty$.

To make these ideas more precise, let us define the *normalized rank* of $Z_0(\theta)$ as

$$\mathcal{R}(\theta) \doteq \frac{1}{m} \left[1 + \sum_{i=1}^{m-1} \mathbb{I}(Z_0(\theta) \prec Z_i(\theta)) \right], \quad (9)$$

where $\mathbb{I}(\cdot)$ is an indicator (its value is 1 if its argument is true and 0 otherwise).

Assume that the target confidence probability can be written as $p = 1 - q/m$ where $0 < q < m$ are user-chosen integers. Then, SPS accepts the null hypothesis, $\theta = \theta^*$, if $\mathcal{R}(\theta) \leq p$, and rejects it if $\mathcal{R}(\theta) > p$. As m and q are free-parameters, they are under our control, hence any (rational) probability can be achieved.

Based on these observations, the *SPS confidence regions* can be defined as

$$\widehat{\Theta}_{n,p} \doteq \left\{ \theta \in \mathbb{R}^d : \mathcal{R}(\theta) \leq p \right\}. \quad (10)$$

It can be proved [2] that these regions have *exact* confidence $\mathbb{P}(\theta^* \in \widehat{\Theta}_{n,p}) = p$. Note that the exact confidence of the regions is guaranteed for *finite samples* despite no knowledge about the particular noise distributions is assumed, moreover, each noise term may have a different distribution with arbitrarily large variance.

There are several important properties of SPS confidence regions [2, 7]. For example, (1) they are *star convex* with the LS estimate as a star center; (2) they are uniformly *strongly consistent*; (3) they have asymptotically the *same size and shape* as the classical confidence ellipsoids; finally (4) they have *ellipsoidal outer approximation* that can be efficiently computed via semidefinite programming problems.

SPS has several generalizations, for example, it can be extended to general stochastic linear (dynamical) systems, even if they are operating in closed-loop [3], and to various non-linear dynamical systems, such as GARCH models [1].

Finally, we note that working with symmetric noises is not crucial for SPS as the theory can be extended to other noise distributions, as long as we know a group of transformations that leave the (joint) distribution of the noises unchanged. For example, one can assume that the noises are *exchangeable* and use random *permutation* matrices as $\{G_i\}$, see [5]. We refer to these generalized variants as (G)SPS.

² Relation “ \prec ” is a total order which we get from “ $<$ ” by random tie-breaking, see [2].

4 Non-Asymptotic Confidence Sets for Regularized Estimates

In this section we are going to extend the theory of (G)SPS, in order to construct non-asymptotic, distribution-free confidence regions around regularized estimates.

First, we consider *ridge regression* (RR) which has the objective function

$$V_{\text{R}}(\boldsymbol{\theta}) \doteq \frac{1}{2} \|y - \boldsymbol{\Phi}\boldsymbol{\theta}\|_2^2 + \frac{\lambda}{2} \|\boldsymbol{\theta}\|_2^2, \quad (11)$$

for a $\lambda \geq 0$ hyper-parameter. It is well-known that RR can be reformulated as LS,

$$\tilde{\boldsymbol{\Phi}} = \begin{bmatrix} \boldsymbol{\Phi} \\ \sqrt{\lambda} I \end{bmatrix}, \quad \text{and} \quad \tilde{y} = \begin{bmatrix} y \\ 0 \end{bmatrix}, \quad (12)$$

where I is the identity matrix, after which we have $V_{\text{R}}(\boldsymbol{\theta}) = 1/2 \|\tilde{y} - \tilde{\boldsymbol{\Phi}}\boldsymbol{\theta}\|_2^2$.

Then, one might be tempted to apply standard SPS to the obtained (ordinary) LS formulation. However, we should proceed with caution, as the new problem has some auxiliary output terms, the zero part of \tilde{y} , to which there are no real noise terms in the original problem. Therefore, the last d terms of the residual vector, $\tilde{y} - \tilde{\boldsymbol{\Phi}}\boldsymbol{\theta}$, should not be perturbed, as the distributional invariance was only assumed for the original noise vector. Consequently, the $\{G_i\}$ matrices should be extended by

$$\tilde{G}_i \doteq \begin{bmatrix} G_i & 0 \\ 0 & I \end{bmatrix}, \quad (13)$$

for $i = 1, \dots, m-1$. Then, using an analogue of (7) to the new LS system with $\{\tilde{G}_i\}$ perturbations, we arrive at the (G)SPS evaluation function for ridge regression,

$$Z_i(\boldsymbol{\theta}) \doteq \left\| \Psi_{\text{R}}^{1/2} [\boldsymbol{\Phi}^T G_i (y - \boldsymbol{\Phi}\boldsymbol{\theta}) - \lambda \boldsymbol{\theta}] \right\|_2^2, \quad (14)$$

where $\Psi_{\text{R}} = (\boldsymbol{\Phi}^T \boldsymbol{\Phi} + \lambda I)^{-1} (\boldsymbol{\Phi}^T \boldsymbol{\Phi}) (\boldsymbol{\Phi}^T \boldsymbol{\Phi} + \lambda I)^{-1}$ is a correction term from the covariance of RR. Based on this evaluation function, *exact* confidence regions can be built around the RR estimate, using the same steps as we had for standard SPS.

Now, let us consider *LASSO* (least absolute shrinkage and selection operator) which applies L1 regularization to enforce *sparsity*. It has the objective function

$$V_{\text{L}}(\boldsymbol{\theta}) \doteq \frac{1}{2} \|y - \boldsymbol{\Phi}\boldsymbol{\theta}\|_2^2 + \lambda \|\boldsymbol{\theta}\|_1, \quad (15)$$

for $\lambda \geq 0$. This objective is no more quadratic and it cannot be traced back to LS. However, the underlying idea of SPS, i.e., to perturb the residuals in the (negative) gradient of the objective, can still be applied. A (sub-) gradient³ of (15) is

$$\nabla_{\boldsymbol{\theta}} V_{\text{L}}(\boldsymbol{\theta}) = \boldsymbol{\Phi}^T \boldsymbol{\Phi} \boldsymbol{\theta} - \boldsymbol{\Phi}^T y + \lambda \text{sign}(\boldsymbol{\theta}), \quad (16)$$

where the sign function is understood component-wise.

³ For our purposes, one of the subgradients is sufficient, thus we do not treat $\nabla_{\boldsymbol{\theta}} V$ set-valued.

Then, we can proceed in the same way as before and perturb the residuals in (16) with $\{G_i\}$, leading to the *(G)SPS evaluation function for LASSO*,

$$Z_i(\theta) \doteq \left\| \Psi_L^{1/2} \left[\Phi^T G_i (y - \Phi\theta) - \lambda \text{sign}(\theta) \right] \right\|_2^2, \quad (17)$$

where Ψ_L is an (optional) correction term, e.g., using the (asymptotic) results of [4], we may use $\Psi_L = (\Phi^T \Phi)^{-1}$. The correction matrix can be interpreted as the square-root of the (estimated) covariance of LASSO (modulo the variance of the noise, as multiplying each Z_i with the same positive scalar does not affect their ordering).

The last method that we discuss is the *elastic net* regularization with objective

$$V_E(\theta) \doteq \frac{1}{2} \|y - \Phi\theta\|_2^2 + \lambda_1 \|\theta\|_1 + \frac{\lambda_2}{2} \|\theta\|_2^2, \quad (18)$$

where $\lambda_1, \lambda_2 \geq 0$ are hyper-parameters. As the objective is the combination of the ridge regression and LASSO objectives, it can be handled using similar ideas. That is we can compute a subgradient of the objective and perturb the residuals based on the transformations $\{G_i\}$ which leave the (joint) distribution of the true noise terms invariant. Then, the *(G)SPS evaluation function for elastic net regularization* is

$$Z_i(\theta) \doteq \left\| \Psi_E^{1/2} \left[\Phi^T G_i (y - \Phi\theta) - \lambda_1 \text{sign}(\theta) - \lambda_2 \theta \right] \right\|_2^2, \quad (19)$$

where Ψ_E can again be an (optional) covariance estimate for the elastic net solution.

The exact confidence of the constructed regions easily follows from the related results for SPS. We leave the investigation of their other properties for further work.

Acknowledgements This research was partially supported by the National Research, Development and Innovation Office (NKFIH), grant numbers ED_18-2-2018-0006 and KH_17_125698, and by the János Bolyai Research Fellowship of the Hung. Academy of Sciences, BO/00217/16/6.

References

1. B. Cs. Csáji. Score permutation based finite sample inference for generalized autoregressive conditional heteroskedasticity (GARCH) models. In *19th International Conference on Artificial Intelligence and Statistics (AISTATS)*, Cadiz, Spain, pages 296–304, 2016.
2. B. Cs. Csáji, M. C. Campi, and E. Weyer. Sign-Perturbed Sums: A new system identification approach for constructing exact non-asymptotic confidence regions in linear regression models. *IEEE Transactions on Signal Processing*, 63(1):169–181, 2015.
3. B. Cs. Csáji and E. Weyer. Closed-loop applicability of the Sign-Perturbed Sums method. In *54th IEEE Conference on Decision and Control, Osaka, Japan*, pages 1441–1446, 2015.
4. K. Knight and W. Fu. Asymptotics for LASSO-type estimators. *Annals of Statistics*, (5):1356–1378, 2000.
5. S. Kolumbán. *System Identification in Highly Non-Informative Environment*. PhD thesis, Budapest Uni. of Techn. and Econ., Hungary, and Vrije Universiteit Brussels, Belgium, 2016.
6. L. Ljung. *System Identification: Theory for the User*. Prentice-Hall, 2nd edition, 1999.
7. E. Weyer, M. C. Campi, and B. Cs. Csáji. Asymptotic properties of SPS confidence regions. *Automatica*, 82:287–294, 2017.