
One-step estimator paths for concave regularization

Matt Taddy

The University of Chicago Booth School of Business

Abstract

The statistics literature of the past 15 years has established many favorable properties for sparse diminishing-bias regularization: techniques which can roughly be understood as providing estimation under penalty functions spanning the range of concavity between L_0 and L_1 norms. However, lasso L_1 -regularized estimation remains the standard tool for industrial ‘Big Data’ applications because of its minimal computational cost and the presence of easy-to-apply rules for penalty selection. In response, this article proposes a simple new algorithm framework that requires no more computation than a lasso path: the path of one-step estimators (POSE) does L_1 penalized regression estimation on a grid of decreasing penalties, but adapts coefficient-specific weights to decrease as a function of the coefficient estimated in the previous path step. This provides sparse diminishing-bias regularization at no extra cost over the fastest lasso algorithms. Moreover, our ‘gamma lasso’ implementation of POSE is accompanied by a reliable heuristic for the fit degrees of freedom, so that standard information criteria can be applied in penalty selection. The methods are illustrated in extensive simulations and in application of logistic regression to evaluating the performance of hockey players.

1 Introduction

For regression in high-dimensions, it is useful to regularize estimation through a penalty on coefficient size. L_1 regularization (i.e., the lasso of Tibshirani, 1996) is especially popular, with costs that are non-differentiable at their minima and can lead to coeffi-

cient solutions of exactly zero. A related approach is concave penalized regularization (e.g. SCAD from Fan and Li 2001 or MCP from Zhang 2010) with cost functions that are also spiked at zero but flatten for large values (as opposed to the constant increase of an L_1 norm). This yields sparse solutions where large non-zero values are estimated with little bias.

The combination of *sparsity* and *diminishing-bias* is appealing in many settings, and a large literature on concave penalized estimation has developed over the past 15 years. For example, many authors (e.g., from Fan and Li 2001 and Fan and Peng 2004) have contributed work on their *oracle properties*, a class of results showing conditions under which coefficient estimates through concave penalization, or in related schemes, will be the same as if you knew the sparse ‘truth’ (either asymptotically or with high probability). From an information compression perspective, the increased sparsity encouraged by diminishing-bias penalties (since single large coefficients are allowed to account for the signals of other correlated covariates) leads to lower memory, storage, and communication requirements. Such savings are very important in distributed computing schemes (e.g., Taddy, 2013).

Unfortunately, exact solvers for concave penalized estimation all require significantly more compute time than a standard lasso. In our experience, this has precluded their use in settings – e.g., in text or web-data analysis – where both n (the number of observations) and p (the number of covariates) are very large. As we review below, recent literature recommends the use of approximate solvers. Certainly, this is necessary for data of the size we encounter in analysis of, say, internet commerce. These approximations take the form of iteratively-weighted- L_1 regularization, where the coefficient-specific weights are based upon estimates of the coefficients taken from previous iterations of the approximate solver. This literature (e.g., Zou and Li, 2008; Fan et al., 2014) holds that even a single step of weighted- L_1 regularization is enough to get solutions that are close to optimal, so long as the pre-estimates are *good enough* starting points. The crux of success with such one-step estimation (OSE) is finding starts that are, indeed, good enough.

This article proposes a framework for sparse

diminishing-bias regularization that combines ideas from OSE with the concept of a *regularization path* – a general technique, most famously associated with the LARS algorithm Efron et al. (2004),

As detailed in Section path of one-step estimators (POSE) does L_1 penalized regression estimation on a grid of decreasing penalties, but adapts coefficient-specific weights to decrease as a function of the coefficient estimated in the previous path step.

2 Path of one-step estimators

Our path of one-step estimators (POSE), in Algorithm 1, uses solutions along the sequence of decreasing penalty sizes, λ^t , as the basis for LLA weights at the next path step. In this, we are assuming a penalty specification such that $\lim_{b \rightarrow 0} c'(|b|) = 1$ and that the cost function is differentiable for $b \neq 0$. This yields a path of one-step LLA penalized coefficient estimates.

Algorithm 1 POSE

Initialize $\hat{\beta}^0 = \mathbf{0}$, so that $\hat{S}_0 = \emptyset$.

Set $\lambda^1 > 0$ with step size $0 < \delta < 1$.

for $t = 1 \dots T$:

$$\omega_j^t = \begin{cases} c'(|\hat{\beta}_j^{t-1}|) & \text{for } j \in \hat{S}_t \\ 1 & \text{for } j \in \hat{S}_t^c \end{cases} \quad (1)$$

$$[\hat{\alpha}, \hat{\beta}]^t = \underset{\alpha, \beta_j \in \mathbb{R}}{\operatorname{argmin}} l(\alpha, \beta) + n \sum_j \lambda^t \omega_j^t |\beta_j| \quad (2)$$

$$\lambda^{t+1} = \delta \lambda^t$$

From an engineering standpoint, POSE has the same appeal as any successful path algorithm: if the estimates change little from iteration t to $t + 1$, then you will be able to quickly solve for a large set of candidate specifications. Following the discussion of Section 3.1, such algorithms are a natural match with one-step estimation: OSE relies upon inputs being close to the optimal solution, which is precisely the setting where path algorithms are most efficient. More rigorously, Theorem ?? applied to POSE yields $\hat{S}_{t-1} \cap S^c = \emptyset \Rightarrow \omega_{S^c}^t = 1$. Thus so long as λ is large enough, Section 3.2 demonstrates that fast diminishing ω_j will help control false discovery and improve prediction. Of course, the moment $\hat{S}_t \cap S^c \neq \emptyset$, diminishing-bias allows spurious covariates to enter with little shrinkage and can move the fit arbitrarily far away from L_0 -optimality – that is, with λ too small the diminishing bias hurts your ability to estimate and predict. This is why it is essential to have a path of candidate λ^t to choose amongst.

3 Sparse regularization paths and diminishing bias

References

- Efron, B., T. Hastie, I. Johnstone, and R. Tibshirani (2004). Least angle regression. *Annals of Statistics* 32, 407–499.
- Fan, J. and R. Li (2001). Variable selection via nonconcave penalized likelihood and its oracle properties. *Journal of the American Statistical Association* 96, 1348–1360.
- Fan, J. and H. Peng (2004). Nonconcave penalized likelihood with a diverging number of parameters. *The Annals of Statistics* 32, 928–961.
- Fan, J., L. Xue, and H. Zou (2014, June). Strong oracle optimality of folded concave penalized estimation. *The Annals of Statistics* 42(3), 819–849.
- Taddy, M. (2013). Distributed multinomial regression. arXiv:1311.6139.
- Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society, Series B* 58, 267–288.
- Zhang, C.-H. (2010). Nearly unbiased variable selection under minimax concave penalty. *Annals of Statistics* 38, 894–942.
- Zou, H. and R. Li (2008). One-step sparse estimates in nonconcave penalized likelihood models. *The Annals of Statistics* 36(4), 1509–1533.