

Model Selection and Estimation in High-dimensional Generalized Linear Models

Francisco Rivera

Jiafeng (Kevin) Chen

December 8, 2018

1 Introduction

The workhorse for high-dimensional regressions is the ℓ_1 lasso penalty ([Tibshirani, 1996](#)).

The original lasso is developed for fitting (normal) linear models with the objective¹

$$\hat{\boldsymbol{\beta}}_{\text{lasso}} = \arg \min_{\boldsymbol{\beta}} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 + \lambda \|\boldsymbol{\beta}\|_1 - \lambda |\beta_0|, \lambda > 0. \quad (1)$$

The first paper about lasso in a GLM setting is [Park and Hastie \(2007\)](#): Consider a scalar GLM with likelihood

$$L(y; \theta, \phi) = \exp \left\{ \frac{y\theta - b(\theta)}{a(\phi)} + c(y, \phi) \right\},$$

where $g(\mathbb{E}[y]) = \eta = \mathbf{x}^\top \boldsymbol{\beta}$ for some scalar function g . Thus we may consider a direct extension of (1):

$$\hat{\boldsymbol{\beta}}_{\text{GLM-lasso}} = \arg \min_{\boldsymbol{\beta}} \underbrace{-\frac{1}{n} \sum_{i=1}^n [y_i \theta(\boldsymbol{\beta})_i - b(\theta(\boldsymbol{\beta})_i)]}_{\ell_n(\boldsymbol{\beta})} + \lambda \|\boldsymbol{\beta}\|_1 - \lambda |\beta_0|, \quad (2)$$

where (2) reduces to (1) if the likelihood is Gaussian and g is the canonical link for the Gaussian model, which is the identity function. Note that the objective (2) is convex if

¹We consider \mathbf{y} an n -vector of response variables whose i th element is y_i . We consider the covariate matrix \mathbf{X} ($n \times p$). In most settings, \mathbf{X} include an intercept column, $\mathbf{x}_0 = \mathbf{1}$, that is not regularized.

we use the canonical link, since the exponential family log-likelihood is concave in θ , and $\theta = \mathbf{x}^\top \boldsymbol{\beta}$ is linear in $\boldsymbol{\beta}$ for the canonical link.

2 Model Selection with Lasso

Brute-force model selection scales poorly with respect to the number of predictors. That is, when building a model to predict \mathbf{y} from a subset of the columns of $\mathbf{X} \in \mathbb{R}^{n,p}$ (since not all columns may be predictive), there are 2^p subsets of columns to regress on. Checking all of these models against each other incurs a cost exponential in p and becomes prohibitively expensive in a high-dimensional setting.

In class, we tackled this problem with forward-selection and backward-deletion, a greedy algorithm that adds and removes covariates until a local optimal is reached. [Park and Hastie \(2007\)](#)'s use of the lasso for GLMs provides a competitive alternative that behaves less greedily. The high-level idea is to note the non-zero coefficients for different penalizations λ and constrain our search to these models. We describe their procedure in this section.

In order to employ this method, we require a method to solve (2). For now, we suppose that we have such a method and can use it as a black box. It is sufficient for us to know that this black box takes in an initial $\boldsymbol{\beta}^{\text{init}}$ and employs a descent-based method to converge toward the optimal $\boldsymbol{\beta}^*$. [Section 3](#) follows-up by describing the inner workings of different approaches.

The lasso penalty induces a sparsity in our optimal $\boldsymbol{\beta}$. This sparsity increases as our penalization term λ becomes bigger. Indeed, if we take $\lambda \rightarrow \infty$, then our solution is driven to $\boldsymbol{\beta} \rightarrow \mathbf{0}$ because the penalization term is all that matters, and the norm of $\boldsymbol{\beta}$ is minimized at $\mathbf{0}$. The entire domain when $\boldsymbol{\beta} = \mathbf{0}$ is uninteresting, so we initialize our algorithm at λ_{max} which we define as the smallest λ such that there is only one non-zero coefficient.

Once we have initialize λ in our algorithm to λ_{max} , we will repeat three steps that [Park and Hastie \(2007\)](#) outline,

1. Determine the step length by which to decrease λ . The intent is to pick this such that the precisely one more coefficient becomes non-zero. Introducing notation, we let λ_k be the λ value we consider in the k th step, and we have to pick $\lambda_k - \lambda_{k+1}$.
2. Using a linear approximation, predict the value of β for the decremented λ_{k+1} .
3. Initialize our solver of (2) with the linear estimate as β^{init} and get the precise solution β^* for the new λ .

doing so until either we have brought λ all the way to 0, or we determine that incremental variables are not improving the model. In order to perform the second step, a linear approximation prescribes that,

$$\beta_{k+1} \approx \beta_k + (\lambda_{k+1} - \lambda_k) \frac{\partial \beta}{\partial \lambda}.$$

Denote by \mathbf{X}_A the matrix made up by the columns of \mathbf{X} whose corresponding β coefficient is non-zero, i.e. in the active set. Furthermore, let \mathbf{W} be the diagonal matrix with i^{th} diagonal element given by,

$$\frac{1}{\mathbb{V}(y_i)} \left(\frac{\partial \mu}{\partial \eta} \right)^2$$

which changes because the derivative is evaluated at a different location based on the value of λ , so we define \mathbf{W}_k as the value of \mathbf{W} corresponding to λ_k . Then, we have that,

$$\beta_{k+1} \approx \beta_k - (\lambda_{k+1} - \lambda_k) (\mathbf{W}_A^\top \mathbf{W}_k \mathbf{X}_A)^{-1} \text{sgn}(\beta_k).$$

Note that in the above equation, we pretend β_k and β_{k+1} only contain their non-zero components, and thus are not going to have all p dimensions. This makes our algebra simpler and we can do this because the linear approximation will not make a previously zero coefficient non-zero.

Finally, we just need to determine our step size. One simple choice could be to use a constant step size. However, in different domains of λ , the same unit change in λ could have dramatically different implications for the sparsity of the resulting $\boldsymbol{\beta}$. If we were working with a lasso in linear regression, \mathbf{W} would not change, and we could pick the step size precisely to coincide with the point at which a coefficient of zero becomes non-zero. [Park and Hastie \(2007\)](#) suggest doing precisely this, with the caveat that because \mathbf{W} does change, it will only be an approximation for the GLM case.

3 Estimation

[Efron et al. \(2004\)](#) gives an efficient estimation procedure for the linear lasso (1) called the Least Angle Regression (LAR), which relies on the fact that the *regularization path*— $\hat{\boldsymbol{\beta}}_i$ as a function of λ —is piecewise linear in (1). Such a structure is often unavailable in applications like (2). The most popular method—proposed by [Friedman, Hastie, and Tibshirani \(2010\)](#) and implemented in R’s `glmnet` package—is *cyclical coordinate descent* with iteratively reweighted least squares. The idea is to approximate $\ell_n(\boldsymbol{\beta})$ in (2) with a second-order Taylor expansion, either globally for all parameters $\boldsymbol{\beta}$ (for scalar-valued GLMs) or locally with a single parameter $\boldsymbol{\beta}_j$ (for vector-valued GLMs, such as the multinomial logistic regression). Such an approximation yields a quadratic function (in the scalar GLM case)

$$\ell_Q(\boldsymbol{\beta}) = \frac{1}{n} \sum_{i=1}^n w_i (y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2.$$

We then solve a local penalized least-squares problem:

$$\boldsymbol{\beta} \leftarrow \arg \min_{\boldsymbol{\beta}} \ell_Q(\boldsymbol{\beta}) + \lambda \|\boldsymbol{\beta}\|_1. \quad (3)$$

via cyclical coordinate descent, i.e. by iteratively solving

$$\beta_j \leftarrow \arg \min_{\beta_j} \ell_Q(\beta) + \lambda \|\beta\|_1, \quad (4)$$

holding all other entries β_{-j} fixed. (4) has an analytical solution for the lasso penalty²

$$\beta_j \leftarrow \frac{S\left(\sum_{i=1}^N w_i x_{ij} (y_i - \tilde{y}_i^{(j)}), \lambda\right)}{\sum_{i=1}^N w_i x_{ij}^2}, \quad S(t, \gamma) = \text{sgn}(t) (|t| - \gamma)_+, \quad \tilde{y}_i^{(j)} = \mathbf{x}_i^\top \beta - x_{ij} \beta_j. \quad (5)$$

We summarize the procedure described above in [Algorithm 1](#).

Algorithm 1 Cyclic coordinate descent algorithm for solving (2) (scalar GLM case) in [Friedman, Hastie, and Tibshirani \(2010\)](#)

```

Initialize  $\beta$ 
for  $\lambda$  on regularization path do
  while  $\beta$  has not converged do
    Approximate  $\ell_n(\beta)$  by  $\ell_Q(\beta)$ 
    while cyclical descent has not converged do
      for  $j$  do
        Update  $\beta_j$  according to (5)
      end for
    end while
  end while
   $\hat{\beta}_\lambda \leftarrow \beta$ 
  Initialize  $\beta$  for next iteration to  $\hat{\beta}_\lambda$ 
end for
    
```

The machine learning literature slightly alters [Algorithm 1](#) and changes (5) into

$$\beta_j \leftarrow S\left(\beta_j - (\nabla_{\beta} \ell_n(\beta))_j \kappa^{-1}, \frac{\lambda}{\kappa}\right)$$

for some *learning rate* $1/\kappa$, in keeping with gradient descent. Moreover, [Shalev-Shwartz and Tewari \(2011\)](#) proves a convergence guarantee for stochastic coordinate descent in this

²[Friedman, Hastie, and Tibshirani \(2010\)](#) show a similar expression for the *elastic net* penalty:

$$\lambda P_\alpha(\beta) = \lambda (\alpha \|\beta\|_1 + (1 - \alpha) \|\beta\|_2).$$

fashion, where, instead of cycling through the coordinates of β , a coordinate is chosen uniformly at random.

Theorem 1. Let $Q(\beta)$ be the objective in (2). At iteration T of the first while-loop in a version of Algorithm 1 with stochastic coordinate descent and gradient updates,

$$\mathbb{E}[Q(\beta_T)] - \mathbb{E}[Q(\hat{\beta}_{\text{GLM-lasso}})] \leq C \frac{p\kappa}{T+1}$$

for constant C a function of the initial starting value $\beta^{(0)}$, assuming that ℓ_n is differentiable with

$$\ell_n(\beta + \eta e_j) \leq \ell_n(\beta) + \eta (\nabla \ell_n)_j + \frac{\kappa}{2} \eta^2$$

for all η, β, j .³

Corollary 2. The runtime to achieve ϵ expected accuracy is bounded by

$$O\left(\frac{np\kappa}{\epsilon} \left\| \hat{\beta}_{\text{GLM-lasso}} \right\|_2^2\right).$$

Moreover, Bradley et al. (2011) show that a parallel version of the coordinate gradient descent procedure above where at each iteration, P (possibly duplicate) coordinates are updated in parallel. For correlated features, such parallelism is dangerous, since updating two correlated features simultaneously may over or undercompensate for the gradient direction. Bradley et al. (2011) quantifies the interference due to correlated features and shows that efficiency increases linearly in the number of parallel processes P so long as $P \leq \frac{p}{\rho}$ where ρ is the largest modulus of the eigenvalues of $X^\top X$.

Coordinate descent methods described above can also become expensive if n is large. The standard machine learning and optimization answer to this problem is to use *stochastic gradient descent*, replacing $\nabla_{\beta} \ell_n(\beta)$ with an unbiased estimate $\mathbf{g}_i = \nabla_{\beta} \log L(y_i; \beta)$, which is

³This condition restricts the choice of κ as a function of the loss criterion.

the gradient evaluated on a single observation.⁴ [Shalev-Shwartz and Tewari \(2011\)](#) consider a mirror descent algorithm in the lasso context, by running stochastic gradient descent on the dual problem and enforcing sparsity in an intelligent manner. Let $\boldsymbol{\gamma} = f(\boldsymbol{\beta})$ be the dual parameter for $\boldsymbol{\beta}$ with an invertible link f . We choose an observation i at random, compute \mathbf{g}_i , and update

$$\begin{aligned}
 \boldsymbol{\gamma} &\leftarrow \boldsymbol{\gamma} - \eta \mathbf{g}_i \\
 \boldsymbol{\gamma}' &\leftarrow \boldsymbol{\gamma} - \eta \lambda \operatorname{sgn}(\boldsymbol{\gamma}) && \text{(Decrease } \|\boldsymbol{\beta}\|_1 \text{)} \\
 \gamma_j &\leftarrow \gamma'_j \mathbb{1}(\operatorname{sgn}(\gamma_j) = \operatorname{sgn}(\gamma'_j)) && \text{(Maintains sparsity)} \\
 \boldsymbol{\beta} &\leftarrow f^{-1}(\boldsymbol{\gamma}).
 \end{aligned}$$

The runtime bound for the stochastic mirror descent algorithm in [Shalev-Shwartz and Tewari \(2011\)](#) is

$$O\left(\frac{p \log p}{\epsilon^2} \left\| \hat{\boldsymbol{\beta}}_{\text{GLM-lasso}} \right\|_2^2\right).$$

We pay the price of the $p \log p$ and ϵ^{-2} dependence, as opposed to p and ϵ^{-1} , in order to achieve the benefit of a n -free runtime.

⁴We can replace this with *batched gradient descent* as well, where the gradient estimate is averaging over a batch of observations.

References

- Bradley, Joseph K, Aapo Kyrola, Danny Bickson, and Carlos Guestrin. 2011. “Parallel coordinate descent for l1-regularized loss minimization.” *arXiv preprint arXiv:1105.5379* .
- Efron, Bradley, Trevor Hastie, Iain Johnstone, Robert Tibshirani et al. 2004. “Least angle regression.” *The Annals of statistics* 32 (2):407–499.
- Friedman, Jerome, Trevor Hastie, and Rob Tibshirani. 2010. “Regularization paths for generalized linear models via coordinate descent.” *Journal of statistical software* 33 (1):1.
- James, Gareth M, Courtney Paulson, and Paat Rusmevichientong. 2013. “Penalized and constrained regression.” *Unpublished Manuscript, available at <http://www-bcf.usc.edu/~gareth/research/Research.html>* .
- Jiang, Wenxin et al. 2007. “Bayesian variable selection for high dimensional generalized linear models: convergence rates of the fitted densities.” *The Annals of Statistics* 35 (4):1487–1511.
- Lockhart, Richard, Jonathan Taylor, Ryan J Tibshirani, and Robert Tibshirani. 2014. “A significance test for the lasso.” *Annals of statistics* 42 (2):413.
- O’sullivan, Finbarr, Brian S Yandell, and William J Raynor Jr. 1986. “Automatic smoothing of regression functions in generalized linear models.” *Journal of the American Statistical Association* 81 (393):96–103.
- Park, Mee Young and Trevor Hastie. 2007. “L1-regularization path algorithm for generalized linear models.” *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 69 (4):659–677.
- Shalev-Shwartz, Shai and Ambuj Tewari. 2011. “Stochastic methods for l1-regularized loss minimization.” *Journal of Machine Learning Research* 12 (Jun):1865–1892.

- Sørensen, Øystein, Kristoffer Herland Hellton, Arnaldo Frigessi, and Magne Thoresen. 2018. “Covariate selection in high-dimensional generalized linear models with measurement error.” *Journal of Computational and Graphical Statistics* (just-accepted).
- Tibshirani, Robert. 1996. “Regression shrinkage and selection via the lasso.” *Journal of the Royal Statistical Society. Series B (Methodological)* :267–288.
- Valdora, Marina, Claudio Agostinelli, and Victor J Yohai. 2017. “Robust Estimation in High Dimensional Generalized Linear Models.” *arXiv preprint arXiv:1709.10261* .
- Van de Geer, Sara, Peter Bühlmann, Yaacov Ritov, Ruben Dezeure et al. 2014. “On asymptotically optimal confidence regions and tests for high-dimensional models.” *The Annals of Statistics* 42 (3):1166–1202.
- Van de Geer, Sara A et al. 2008. “High-dimensional generalized linear models and the lasso.” *The Annals of Statistics* 36 (2):614–645.
- Xu, Jason, Eric Chi, and Kenneth Lange. 2017. “Generalized Linear Model Regression under Distance-to-set Penalties.” In *Advances in Neural Information Processing Systems*. 1385–1395.