
Supplemental Material: Generalized Linear Model Regression under Distance-to-set Penalties

Jason Xu

University of California, Los Angeles
jqxu@ucla.edu

Eric C. Chi

North Carolina State University
eric_chi@ncsu.edu

Kenneth Lange

University of California, Los Angeles
klange@ucla.edu

1 Proof of Convergence

We repeat the statement of Theorem 3.1 below:

Theorem 1.1. *Consider the algorithm map*

$$\mathcal{M}(\beta) = \beta - \eta_\beta \mathbf{H}(\beta)^{-1} \nabla f(\beta),$$

where the step size η_β has been selected by Armijo backtracking. Assume that f is coercive: $\lim_{\|\beta\| \rightarrow \infty} f(\beta) = +\infty$. Then the limit points of the sequence $\beta_{k+1} = \mathcal{M}(\beta_k)$ are stationary points of $f(\beta)$. Moreover, this set of limit points is compact and connected.

Our algorithm selects the step-size η according to the Armijo condition: suppose \mathbf{v} is a descent direction at β in the sense that $df(\beta)\mathbf{v} < 0$. The Armijo condition chooses a step size η such that

$$f(\beta + t\mathbf{v}) \leq f(\beta) + \alpha \eta df(\beta)\mathbf{v},$$

for a constant $\alpha \in (0, 1)$. Before proving the statement, the next lemma follows an argument in Chapter 12 of [3] to show that step-halving under the Armijo condition requires finitely many steps. We may then apply a similar argument used in [6].

Lemma 1.2. *Given $\alpha \in (0, 1)$ and $\sigma \in (0, 1)$, there exists an integer $s \geq 0$ such that*

$$f(\beta + \sigma^s \mathbf{v}) \leq f(\beta) + \alpha \sigma^s df(\beta)\mathbf{v},$$

where $\mathbf{v} = -\mathbf{H}(\beta)^{-1} \nabla f(\beta)$.

Proof. Since f is coercive by assumption, its sublevel sets are compact. Namely, the set $\mathcal{S}_f(\beta_0) \equiv \{\beta : f(\beta) \leq f(\beta_0)\}$ is compact. Smoothness of the GLM likelihood and squared distance penalty ensure continuity of $\nabla f(\beta)$ and $\mathbf{H}(\beta)$. Together with coercivity, this implies that there exist positive constants a and b , such that

$$\|\mathbf{H}(\beta)^{-1}\| \leq a; \quad \|\mathbf{H}(\beta)\| \leq b; \quad \|d^2 \mathcal{L}(x)\| \leq c$$

for all $\beta \in \mathcal{S}_f(\beta_0)$. Together with the fact that the Euclidean distance to a closed set $\text{dist}(\beta, C)$ is a Lipschitz function with Lipschitz constant 1, we produce the inequality

$$f(\beta + \eta \mathbf{v}) \leq f(\beta) + \eta df(\beta)\mathbf{v} + \frac{1}{2} \eta^2 L \|\mathbf{v}\|^2 \tag{1}$$

where $L = 1 + c$. The squared term appearing at the end of (1) can be bounded by

$$\|\mathbf{v}\|^2 = \|\mathbf{H}(\beta)^{-1} \nabla f(\beta)\|^2 \leq a^2 \|\nabla f(\beta)\|^2.$$

We next identify a bound for $\|\nabla f(\beta)\|^2$:

$$\begin{aligned}
\|\nabla f(\beta)\|^2 &= \|\mathbf{H}(\beta)^{1/2} \mathbf{H}(\beta)^{-1/2} \nabla f(\beta)\|^2 \\
&\leq \|\mathbf{H}(\beta)^{1/2}\|^2 \|\mathbf{H}(\beta)^{-1/2} \nabla f(\beta)\|^2 \\
&\leq bdf(\beta) \mathbf{H}(\beta)^{-1} \nabla f(\beta) \\
&= -bdf(\beta) \mathbf{v}.
\end{aligned} \tag{2}$$

Combining inequalities (1) and (2) yields

$$f(\beta + \eta \mathbf{v}) \leq f(\beta) + \eta \left(1 - \frac{a^2 b L}{2} t\right) df(\beta) \mathbf{v}.$$

Thus, the Armijo condition is guaranteed to be satisfied as soon as $s \geq s^*$, where

$$s^* = \frac{1}{\ln \sigma} \ln \left(\frac{2(1 - \alpha)}{a^2 b L} \right).$$

Of course, in practice a much lower value of s may suffice. □

We are now ready to prove the original theorem: note the argument applies whenever C is convex.

Proof. Consider the iterates of the algorithm $\beta_{k+1} = \mathcal{M}(\beta_k) = \beta_k + \sigma^{s_k} \mathbf{v}_k$. Since $f(\beta)$ is continuous, f attains its infimum over $\mathcal{S}_f(\beta_0)$, and therefore the monotonically decreasing sequence $f(\beta_k)$ is bounded below. This implies that $f(\beta_k) - f(\beta_{k+1})$ converges to 0. Let s_k denote the number of backtracking steps taken at the k th iteration under the Armijo stopping rule. By Lemma 1.2, s_k is finite, and thus

$$\begin{aligned}
f(\beta_k) - f(\beta_{k+1}) &\geq -\alpha \sigma^{s_k} df(\beta_k) \mathbf{v}_k \\
&= \alpha \sigma^{s_k} df(\beta_k) \mathbf{H}(\beta_k)^{-1} \nabla f(\beta_k) \\
&\geq \frac{\alpha \sigma^{s_k}}{\beta} \|\nabla f(\beta_k)\|^2 \\
&\geq \frac{\alpha \sigma^{s^*+1}}{\beta} \|\nabla f(\beta_k)\|^2.
\end{aligned}$$

This inequality implies that $\|\nabla f(\beta_k)\|$ converges to 0, and therefore all the limit points of the sequence β_k are stationary points of $f(\beta)$. Further, taking norms of the update yields the inequality

$$\begin{aligned}
\|\beta_{k+1} - \beta_k\| &= \sigma^{s_k} \|\mathbf{H}(\beta_k)^{-1} \nabla f(\beta_k)\| \\
&\leq \sigma^{s_k} a \|\nabla f(\beta_k)\|.
\end{aligned}$$

Thus, the iterates β_k are a bounded sequence such that $\|\beta_{k+1} - \beta_k\|$ tends to 0, allowing us to conclude that the limit points form a compact and connected set by Propositions 12.4.2 and 12.4.3 in [3]. □

2 Bregman Divergences

Let $\phi : \Omega \mapsto \mathbb{R}$ be a strictly convex function defined on a convex domain $\Omega \subset \mathbb{R}^n$ differentiable on the interior of Ω . The Bregman divergence [1] between \mathbf{u} and \mathbf{v} with respect to ϕ is defined as

$$D_\phi(\mathbf{v}, \mathbf{u}) = \phi(\mathbf{v}) - \phi(\mathbf{u}) - d\phi(\mathbf{u})(\mathbf{v} - \mathbf{u}). \tag{3}$$

Note that the Bregman divergence (3) is a convex function of its first argument \mathbf{v} , and measures the distance between \mathbf{v} and a first order Taylor expansion of ϕ about \mathbf{u} evaluated at \mathbf{v} . While the Bregman divergence is not a metric as it is not symmetric in general, it provides a natural notion of directed distance. It is non-negative for all \mathbf{u}, \mathbf{v} and equal to zero if and only if $\mathbf{v} = \mathbf{u}$. Instances of Bregman divergences abound in statistics and machine learning, many useful measures of closeness.

Recall exponential family distributions takes the canonical form

$$p(y|\theta, \tau) = C_1(y, \tau) \exp \left\{ \frac{y\theta - \psi(\theta)}{C_2(\tau)} \right\}.$$

Each distribution belonging to an exponential family shares a close relationship to a Bregman divergence, and we may explicitly relate GLMs as a special case using this connection. Specifically, the conjugate of its cumulant function ψ , which we denote ζ , uniquely generates a Bregman divergence D_ζ that represents the exponential family likelihood up to proportionality [5]. With g denoting the link function, the negative log-likelihood of y can be written as its Bregman divergence to the mean:

$$-\ln p(y|\theta, \tau) = D_\zeta(y, g^{-1}(\theta)) + C(y, \tau).$$

As an example, the cumulant function in the Poisson likelihood is $\psi(x) = e^x$, whose conjugate $\zeta(x) = x \ln x - x$ produces the relative entropy

$$D_\zeta(p, q) = p \ln(p/q) - p + q.$$

Similarly, recall that the Bernoulli likelihood in logistic regression has cumulant function $\psi(x) = \ln(1 + \exp(x))$. Its conjugate is given by $\zeta(x) = x \ln x + (1 - x) \ln(1 - x)$, and generates

$$D_\zeta(p, q) = p \ln \frac{p}{q} + (1 - p) \ln \frac{1 - p}{1 - q}.$$

This relationship implies that maximizing the likelihood in an exponential family is equivalent to minimizing a corresponding Bregman divergence between the data \mathbf{y} and the regression coefficients β . Notice this is a different statement than the well-known equivalence between maximizing the likelihood and minimizing the Kullback-Leibler divergence between the empirical and parametrized distributions. The gradients of the Bregman projection take the form

$$\nabla D_\phi(\mathcal{P}_{C_i}^\phi(\beta_k), \beta) = d^2\phi(\beta) (\beta - \mathcal{P}_{C_i}^\phi(\beta_k)).$$

Further, the notion of Bregman divergence naturally applies to matrices:

$$D_\phi(\mathbf{V}, \mathbf{U}) = \phi(\mathbf{V}) - \phi(\mathbf{U}) - \langle \nabla \phi(\mathbf{U}), \mathbf{V} - \mathbf{U} \rangle$$

where $\langle \mathbf{V}, \mathbf{U} \rangle = \text{Tr}(\mathbf{V}\mathbf{U}^T)$ denotes the inner product. For instance, the squared Frobenius distance between \mathbf{V}, \mathbf{U} is generated by the choice of $\phi(\mathbf{V}) = \frac{1}{2} \|\mathbf{V}\|_F^2$. The MM algorithm therefore applies analogously to objective functions consisting of multiple Bregman divergences.

3 EEG dataset

The dataset we consider using rank restricted matrix regression seeks to study the association between alcoholism and the voltage patterns over times and channels from EEG data. The data are collected by [7], who provide further details of the experiment, and measures subjects over 120 trials. The study consists of 77 individuals with alcoholism and 45 controls. For each subject, 64 channels of electrodes were placed across the scalp, and voltages are recorded at 256 time points sampled at 256 Hz over one second. This is repeated over 120 trials with three different stimuli. Following the practice of previous studies of the data by [4, 2, 8], we consider covariates \mathbf{X} representing the average over all trials of voltages recorded from each electrode. Other than averaging over trials, no data preprocessing is applied. \mathbf{X} is thus a 256×64 matrix whose ij th entries represent the voltage at time i in channel or electrode j , averaged over the 120 trials. The binary responses y_i indicate whether subject i has alcoholism.

As mentioned in the main text, the study by [4] focuses on reduction of the data via dimension folding, and the matrix-variate logistic regression algorithm proposed by [2] is also applied to preprocessed data using a generic dimension reduction technique. The nuclear norm shrinkage proposed by [8] is the first to consider matrix regression on the full, unprocessed data (apart from averaging over the 120 trials). The authors [8] point out that previous methods nonetheless attain better classification rates, likely due to the fact that preprocessing and tuning were chosen to optimize predictive accuracy. Indeed, the lowest misclassification rate reported in previous analyses is 0.139 by [2], yet the authors show that their method is equivalent to seeking the best rank 1 approximation to the true coefficient matrix in terms of Kullback-Leibler divergence. Since this approach is strictly more restrictive than ours, which attains an error of 0.1475, we agree with [8] in concluding that the lower misclassification error achieved by previous studies can be largely attributed to benefiting from removal of noise via data preprocessing and dimension reduction.

4 Additional comparisons in Gaussian regression

We consider an analogous simulation study including a comparison to the two-stage relaxed LASSO procedure, implemented in R package `relaxo`. The author’s implementation is limited to the Gaussian case, and we consider linear regression with dimension $n = 2000$ as the number of samples m varies, with $k = 12$ nonzero true coefficients. We consider a reduced experiment due to runtime considerations of `relaxo`, repeating only over 20 trials and varying m by increments of 200. Though timing is heavily dependent on implementations, the average total runtimes (across all values of m) of the experiment across trials for MM, MCP, SCAD, and relaxed lasso are 96.8, 137.5, 107.3, 4876.6 seconds, respectively. We see that relaxed LASSO is effective toward removing the bias induced by standard LASSO, and overall results are similar to those included in the main text.

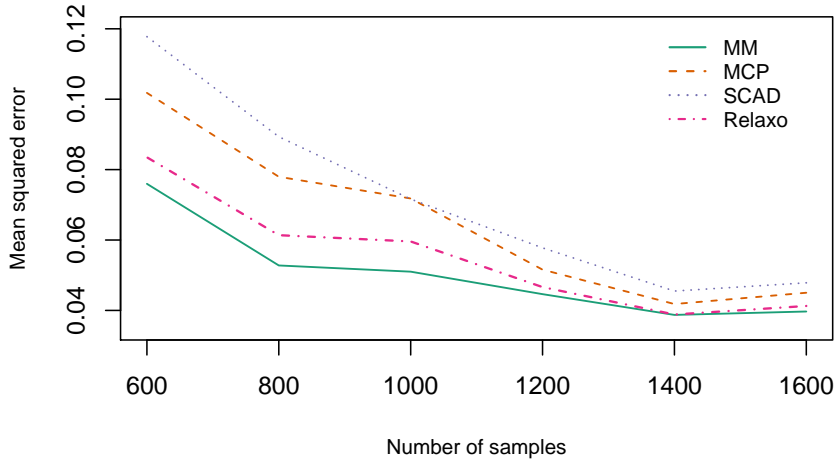


Figure 1: Median MSE over 20 trials as a function of the number of samples m in linear regression under our MM approach, the two-stage relaxed LASSO procedure, SCAD and MCP.

References

- [1] Bregman, L. M. The relaxation method of finding the common point of convex sets and its application to the solution of problems in convex programming. *USSR Computational Mathematics and Mathematical Physics*, 7(3):200–217, 1967.
- [2] Hung, H. and Wang, C.-C. Matrix variate logistic regression model with application to EEG data. *Biostatistics*, 14(1):189–202, 2013.
- [3] Lange, K. *Optimization*. Springer Texts in Statistics. Springer-Verlag, New York, 2nd edition, 2013.
- [4] Li, B., Kim, M. K., and Altman, N. On dimension folding of matrix-or array-valued statistical objects. *The Annals of Statistics*, pages 1094–1121, 2010.
- [5] Polson, N. G., Scott, J. G., and Willard, B. T. Proximal algorithms in statistics and machine learning. *Statistical Science*, 30(4):559–581, 2015.
- [6] Xu, J., Chi, E. C., Yang, M., and Lange, K. A majorization-minimization algorithm for split feasibility problems. *arXiv preprint arXiv:1612.05614*, 2017.
- [7] Zhang, X. L., Begleiter, H., Porjesz, B., Wang, W., and Litke, A. Event related potentials during object recognition tasks. *Brain Research Bulletin*, 38(6):531–538, 1995.
- [8] Zhou, H. and Li, L. Regularized matrix regression. *Journal of the Royal Statistical Society: Series B (Methodological)*, 76(2):463–483, 2014.