

Optimizing Approximate Leave-one-out Cross-validation to Tune Hyperparameters

Ryan Burn

RYAN.BURN@GMAIL.COM

*Department of Mathematics
University of Washington
Seattle, WA 98195-4350, USA*

Editor:

Abstract

For a large class of regularized models, leave-one-out cross-validation (LO) can be efficiently estimated with an approximate leave-one-out formula (ALO). We consider the problem of adjusting hyperparameters so as to maximize performance on ALO. We derive efficient formulas to compute the gradient and hessian of ALO and show how to apply a second-order optimizer to find hyperparameters. We demonstrate the usefulness of the approach by finding hyperparameters for regularized logistic regression and ridge regression on a variety of different real-world data sets.

Keywords: Hyperparameter Optimization, Logistic Regression, Ridge Regression, Machine Learning, Python

1. Introduction

Let $\mathcal{D} = \{(y_1, \mathbf{x}_1), \dots, (y_n, \mathbf{x}_n)\}$ denote a data set where $\mathbf{x}_i \in \mathbb{R}^p$ are features and $y_i \in \mathbb{R}$ are responses. In many applications, we model the observations as independent identically distributed draws from an unknown joint distribution of the form $q(y_i, \mathbf{x}_i) = q_1(y_i | \mathbf{x}_i^\top \boldsymbol{\beta}_*) q_2(\boldsymbol{\beta}_*)$, and we estimate $\boldsymbol{\beta}_*$ using the optimization problem

$$\hat{\boldsymbol{\beta}} \triangleq \underset{\boldsymbol{\beta}}{\operatorname{argmin}} \left\{ \sum_{i=1}^n \ell(y_i | \mathbf{x}_i^\top \boldsymbol{\beta}) + R_{\boldsymbol{\lambda}}(\boldsymbol{\beta}) \right\} \quad (1)$$

where ℓ is a loss function, $R_{\boldsymbol{\lambda}}(\boldsymbol{\beta}) = \sum_{j=1}^p r_{j,\boldsymbol{\lambda}}(\hat{\beta}_j)$ is a regularizer, and $\boldsymbol{\lambda}$ are hyperparameters. $R_{\boldsymbol{\lambda}}$ controls the complexity of the model by penalizing larger values of $\boldsymbol{\beta}$ and can prevent overfitting. We aim to select $\boldsymbol{\lambda}$ so as to minimize the out-of-sample prediction error

$$\operatorname{Err}_{\text{out}} \triangleq \mathbb{E} \left[\ell(y_o | \mathbf{x}_o^\top \hat{\boldsymbol{\beta}}) | \mathcal{D} \right]$$

where (y_o, \mathbf{x}_o) is an unseen sample from the distribution $q(y, \mathbf{x})$. Because q is unknown, we estimate $\operatorname{Err}_{\text{out}}$ with a function $f(\boldsymbol{\lambda} | \mathcal{D})$ and apply a hyperparameter optimization algorithm to find

$$\hat{\boldsymbol{\lambda}} \triangleq \underset{\boldsymbol{\lambda}}{\operatorname{argmin}} f(\boldsymbol{\lambda} | \mathcal{D})$$

Empirical evidence has shown leave-one-out cross-validation (LO) to be an accurate method for estimating Err_{out} (Rad et al., 2020). In the general case, LO can be expensive to compute, requiring a model to be fit n times; but under certain conditions, it can be efficiently estimated with a closed-form approximate leave-one-out formula (ALO), and the approximation has been shown to be accurate in high-dimensional settings (Rad and Maleki, 2020). Additionally, for the special case of ridge regression, ALO is exact and equivalent to Allen’s PRESS (Allen, 1974).

In this paper, we derive efficient formulas to compute the gradient and hessian of ALO given a smooth loss function and a smooth regularizer, and we show how to apply trust-region optimization to find hyperparameters for several different classes of regularized models.

1.1 Relevant Work

Grid search is a commonly used approach for hyperparameter optimization (HO). While it can work well for models with only a single hyperparameter, it requires a search space to be specified and quickly becomes inefficient when multiple hyperparameters are used (Bergstra and Bengio, 2012). Other scalable gradient-based approaches have focused on using holdout sets to approximate Err_{out} (Do et al., 2008; Bengio, 2000). By using ALO, we expect our objective to be a more accurate proxy for Err_{out} . We further improve on previous approaches by computing the hessian of our objective function and applying trust-region optimization, allowing us to make global convergence guarantees.

1.2 Notation

We denote vectors by lowercase bold letters and matrices by uppercase bold letters. We use the notation $\text{vec}[\{a_i\}_i]$ to denote the column vector $(a_1, a_2, \dots)^\top$. For a function $f : \mathbb{R} \rightarrow \mathbb{R}$, we denote its 1st, 2nd, 3rd, and 4th derivatives by \dot{f} , \ddot{f} , \dddot{f} , and \ddddot{f} , respectively. Finally, given a function $f : \mathbb{R}^p \rightarrow \mathbb{R}$, we denote the gradient by ∇f and the hessian by $\nabla^2 f$.

2. Preliminary: Trust-region Optimization

Let $f : \mathbb{R}^p \rightarrow \mathbb{R}$ denote a twice-differentiable objective function. Trust-region optimization is an iterative, second-order optimization algorithm that produces a sequence $\{\mathbf{x}_k\}$ where the k^{th} iteration is generated by updating the previous iteration with a solution to the subproblem (Sorensen, 1982)

$$\begin{aligned} \mathbf{x}_k &= \mathbf{x}_{k-1} + \hat{\mathbf{s}}_k \\ \hat{\mathbf{s}}_k &= \underset{\mathbf{s}}{\text{argmin}} \left\{ f(\mathbf{x}_{k-1}) + \nabla f(\mathbf{x}_{k-1})^\top \mathbf{s} + \frac{1}{2} \mathbf{s}^\top \nabla^2 f(\mathbf{x}_{k-1}) \mathbf{s} \right\} \\ \text{s.t. } & \|\mathbf{s}\| \leq \delta_k \end{aligned}$$

The subproblem minimizes the second-order approximation of f at \mathbf{x}_{k-1} within the neighborhood $\|\mathbf{s}\| \leq \delta_k$, called the trust-region. With the trust-region, we can restrict the second-order approximation to areas where it models f well. Efficient algorithms exist to solve the subproblem regardless of whether $\nabla^2 f(\mathbf{x}_{k-1})$ is positive-definite, making trust-region optimization well-suited for non-convex optimization problems (Moré and Sorensen,

1983). With proper rules for updating δ_k and standard assumptions, such as Lipschitz continuity of ∇f , trust-region optimization is globally convergent. Additionally, if $\nabla^2 f$ is Lipschitz continuous for all \mathbf{x} sufficiently close to a nondegenerate second-order stationary point \mathbf{x}_* where $\nabla^2 f(\mathbf{x}_*)$ is positive-definite, then trust-region optimization has quadratic local convergence (Nocedal and Wright, 1999).

3. Preliminary: The ALO Formula

Put $\ell_i(u) \triangleq \ell(y_i | u)$ and $u_i = \mathbf{x}_i^\top \hat{\boldsymbol{\beta}}$. The LO estimate for Err_{out} is defined as

$$\text{LO}_{\boldsymbol{\lambda}} \triangleq \frac{1}{n} \sum_{i=1}^n \ell_i(\mathbf{x}_i^\top \hat{\boldsymbol{\beta}}_{/i})$$

where

$$\hat{\boldsymbol{\beta}}_{/i} \triangleq \underset{\boldsymbol{\beta}}{\text{argmin}} \left\{ \sum_{j \neq i} \ell_j(\mathbf{x}_j^\top \boldsymbol{\beta}) + R_{\boldsymbol{\lambda}}(\boldsymbol{\beta}) \right\} \quad (2)$$

ALO works by finding $\hat{\boldsymbol{\beta}}$ and then using Newton's method to approximate $\hat{\boldsymbol{\beta}}_{/i}$ from the gradient and hessian of Equation 2 at $\hat{\boldsymbol{\beta}}$. Let $\mathbf{g}_{/i}$ and $\mathbf{H}_{/i}$ denote the gradient and hessian of Equation 2 at $\hat{\boldsymbol{\beta}}$. The ALO approximation to $\hat{\boldsymbol{\beta}}_{/i}$ is $\tilde{\boldsymbol{\beta}}_{/i} \triangleq \hat{\boldsymbol{\beta}} - \mathbf{H}_{/i}^{-1} \mathbf{g}_{/i}$ and

$$\text{ALO}_{\boldsymbol{\lambda}} \triangleq \sum_{i=1}^n \ell_i(u_i - \mathbf{x}_i^\top \mathbf{H}_{/i}^{-1} \mathbf{g}_{/i})$$

Computed naively, this formula would require solving a linear system of order p n times, but applying the matrix inversion lemma leads to a more efficient form. Let \mathbf{X} represent the feature matrix whose i^{th} row is \mathbf{x}_i ; let \mathbf{H} denote the hessian of Equation 1 at $\hat{\boldsymbol{\beta}}$. Define $\mathbf{W} \triangleq \nabla^2 R_{\boldsymbol{\lambda}}(\hat{\boldsymbol{\beta}})$. Then

$$\mathbf{H} = \mathbf{X}^\top \mathbf{A} \mathbf{X} + \mathbf{W}$$

where \mathbf{A} is a diagonal matrix with $\mathbf{A}_{ii} = \ddot{\ell}_i(u_i)$ and

$$\text{ALO}_{\boldsymbol{\lambda}} = \frac{1}{n} \sum_{i=1}^n \ell_i \left(u_i + \frac{\dot{\ell}_i(u_i) h_i}{1 - \ddot{\ell}_i(u_i) h_i} \right)$$

where $h_i \triangleq \mathbf{x}_i^\top \mathbf{H}^{-1} \mathbf{x}_i$. See Rad and Maleki (2020) for a derivation.

4. Optimizing ALO

Put

$$\begin{aligned} \tilde{u}_i(u, h) &\triangleq u + \frac{\dot{\ell}_i(u) h}{1 - \ddot{\ell}_i(u) h} \\ \tilde{\ell}_i(u, h) &\triangleq \ell_i(\tilde{u}_i(u, h)) \\ f_i(\boldsymbol{\lambda}) &\triangleq \tilde{\ell}_i(u_i, h_i) \\ f(\boldsymbol{\lambda}) &\triangleq \sum_{i=1}^n f_i(\boldsymbol{\lambda}) \end{aligned}$$

Minimizing f is equivalent to minimizing ALO_{λ} . And the gradient and hessian of f can be computed by summing over the gradient and hessian of f_i .

$$\begin{aligned}\frac{\partial f}{\partial \lambda_s} &= \sum_{i=1}^n \frac{\partial f_i}{\partial \lambda_s} \\ \frac{\partial^2 f}{\partial \lambda_s \partial \lambda_t} &= \sum_{i=1}^n \frac{\partial^2 f_i}{\partial \lambda_s \partial \lambda_t}\end{aligned}$$

We present formulas for computing ∇f_i and $\nabla^2 f_i$. See Appendix A and Appendix B for derivations. In general, $\nabla^2 f$ will not be positive-definite; but by applying a trust-region optimization algorithm, we can find local minimums.

Theorem 1 *The gradient of f_i can be computed as*

$$\frac{\partial f_i}{\partial \lambda_s} = \frac{\partial \tilde{\ell}_i}{\partial u}(u_i, h_i) \times \frac{\partial u_i}{\partial \lambda_s} + \frac{\partial \tilde{\ell}_i}{\partial h}(u_i, h_i) \times \frac{\partial h_i}{\partial \lambda_s}$$

where

$$\begin{aligned}\frac{\partial u_i}{\partial \lambda_s} &= \mathbf{x}_i^\top \frac{\partial \hat{\beta}}{\partial \lambda_s} \\ \frac{\partial \hat{\beta}}{\partial \lambda_s} &= -\mathbf{H}^{-1} \frac{\partial \nabla R_{\lambda}}{\partial \lambda_s}(\hat{\beta}) \\ \frac{\partial h_i}{\partial \lambda_s} &= -\mathbf{x}_i^\top \mathbf{H}^{-1} \frac{\partial \mathbf{H}}{\partial \lambda_s} \mathbf{H}^{-1} \mathbf{x}_i \\ \frac{\partial \mathbf{H}}{\partial \lambda_s} &= \mathbf{X}^\top \frac{\partial \mathbf{A}}{\partial \lambda_s} \mathbf{X} + \frac{\partial \mathbf{W}}{\partial \lambda_s}\end{aligned}$$

The diagonal matrices \mathbf{A} and \mathbf{W} have derivatives

$$\begin{aligned}\left(\frac{\partial \mathbf{A}}{\partial \lambda_s}\right)_{ii} &= \ddot{\ell}_i(u_i) \times \frac{\partial u_i}{\partial \lambda_s} \\ \left(\frac{\partial \mathbf{W}}{\partial \lambda_s}\right)_{jj} &= \frac{\partial \ddot{r}_{j,\lambda}}{\partial \lambda_s}(\hat{\beta}_j) + \ddot{r}_{j,\lambda}(\hat{\beta}_j) \times \frac{\partial \hat{\beta}_j}{\partial \lambda_s}\end{aligned}$$

Finally, $\tilde{\ell}_i$ has derivatives

$$\begin{aligned}\frac{\partial \tilde{\ell}_i}{\partial u} &= \dot{\ell}_i(\tilde{u}_i(u, h)) \times \frac{\partial \tilde{u}_i}{\partial u} \\ \frac{\partial \tilde{\ell}_i}{\partial h} &= \dot{\ell}_i(\tilde{u}_i(u, h)) \times \frac{\partial \tilde{u}_i}{\partial h}\end{aligned}$$

and \tilde{u}_i has derivatives

$$\begin{aligned}\frac{\partial \tilde{u}_i}{\partial u} &= \frac{1}{1 - \ddot{\ell}_i(u) h} + \frac{\dot{\ell}_i(u) \ddot{\ell}_i(u) h^2}{(1 - \ddot{\ell}_i(u) h)^2} \\ \frac{\partial \tilde{u}_i}{\partial h} &= \frac{\dot{\ell}_i(u)}{(1 - \ddot{\ell}_i(u) h)^2}\end{aligned}$$

Theorem 2 *The hessian of f_i can be computed as*

$$\begin{aligned} \frac{\partial^2 f_i}{\partial \lambda_s \partial \lambda_t} &= \frac{\partial \tilde{\ell}_i}{\partial u}(u_i, h_i) \times \frac{\partial^2 u_i}{\partial \lambda_s \partial \lambda_t} + \frac{\partial^2 \tilde{\ell}_i}{\partial u^2}(u_i, h_i) \times \frac{\partial u_i}{\partial \lambda_s} \times \frac{\partial u_i}{\partial \lambda_t} \\ &\quad + \frac{\partial^2 \tilde{\ell}_i}{\partial u \partial h}(u_i, h_i) \times \left[\frac{\partial u_i}{\partial \lambda_s} \times \frac{\partial h_i}{\partial \lambda_t} + \frac{\partial u_i}{\partial \lambda_t} \times \frac{\partial h_i}{\partial \lambda_s} \right] \\ &\quad + \frac{\partial \tilde{\ell}_i}{\partial h}(u_i, h_i) \times \frac{\partial^2 h_i}{\partial \lambda_s \partial \lambda_t} + \frac{\partial^2 \tilde{\ell}_i}{\partial h^2}(u_i, h_i) \times \frac{\partial h_i}{\partial \lambda_s} \times \frac{\partial h_i}{\partial \lambda_t} \end{aligned}$$

where

$$\begin{aligned} \frac{\partial^2 u_i}{\partial \lambda_s \partial \lambda_t} &= \mathbf{x}_i^\top \frac{\partial^2 \hat{\beta}}{\partial \lambda_s \partial \lambda_t} \\ \frac{\partial^2 \hat{\beta}}{\partial \lambda_s \partial \lambda_t} &= -\mathbf{H}^{-1} \mathbf{X}^\top \cdot \text{vec} \left[\left\{ \ddot{\ell}_i(u_i) \times \frac{\partial u_i}{\partial \lambda_s} \times \frac{\partial u_i}{\partial \lambda_t} \right\}_i \right] \\ &\quad - \mathbf{H}^{-1} \frac{\partial \nabla^2 R_\lambda}{\partial \lambda_s}(\hat{\beta}) \frac{\partial \hat{\beta}}{\partial \lambda_t} - \mathbf{H}^{-1} \frac{\partial \nabla^2 R_\lambda}{\partial \lambda_t}(\hat{\beta}) \frac{\partial \hat{\beta}}{\partial \lambda_s} - \mathbf{H}^{-1} \frac{\partial \nabla R_\lambda}{\partial \lambda_s \partial \lambda_t}(\hat{\beta}) \\ &\quad - \mathbf{H}^{-1} \cdot \text{vec} \left[\left\{ \ddot{r}_{j,\lambda}(\hat{\beta}_j) \times \frac{\partial \hat{\beta}_j}{\partial \lambda_s} \times \frac{\partial \hat{\beta}_j}{\partial \lambda_t} \right\}_j \right] \\ \frac{\partial^2 h_i}{\partial \lambda_s \partial \lambda_t} &= 2\mathbf{x}_i^\top \mathbf{H}^{-1} \frac{\partial \mathbf{H}}{\partial \lambda_s} \mathbf{H}^{-1} \frac{\partial \mathbf{H}}{\partial \lambda_t} \mathbf{H}^{-1} \mathbf{x}_i - \mathbf{x}_i^\top \mathbf{H}^{-1} \frac{\partial^2 \mathbf{H}}{\partial \lambda_s \partial \lambda_t} \mathbf{H}^{-1} \mathbf{x}_i \\ \frac{\partial^2 \mathbf{H}}{\partial \lambda_s \partial \lambda_t} &= \mathbf{X}^\top \frac{\partial^2 \mathbf{A}}{\partial \lambda_s \partial \lambda_t} \mathbf{X} + \frac{\partial^2 \mathbf{W}}{\partial \lambda_s \partial \lambda_t} \end{aligned}$$

The diagonal matrices \mathbf{A} and \mathbf{W} have second derivatives

$$\begin{aligned} \left(\frac{\partial^2 \mathbf{A}}{\partial \lambda_s \partial \lambda_t} \right)_{ii} &= \ddot{\ell}_i(u_i) \times \frac{\partial^2 u_i}{\partial \lambda_s \partial \lambda_t} + \ddot{\ell}_i(u_i) \times \frac{\partial u_i}{\partial \lambda_s} \times \frac{\partial u_i}{\partial \lambda_t} \\ \left(\frac{\partial^2 \mathbf{W}}{\partial \lambda_s \partial \lambda_t} \right)_{jj} &= \frac{\partial^2 \ddot{r}_{j,\lambda}}{\partial \lambda_s \partial \lambda_t}(\hat{\beta}_j) + \frac{\partial \ddot{r}_{j,\lambda}}{\partial \lambda_s}(\hat{\beta}_j) \times \frac{\partial \hat{\beta}_j}{\partial \lambda_t} + \frac{\partial \ddot{r}_{j,\lambda}}{\partial \lambda_t}(\hat{\beta}_j) \times \frac{\partial \hat{\beta}_j}{\partial \lambda_s} \\ &\quad + \ddot{r}_{j,\lambda}(\hat{\beta}_j) \times \frac{\partial^2 \hat{\beta}_j}{\partial \lambda_s \partial \lambda_t} + \ddot{r}_{j,\lambda}(\hat{\beta}_j) \times \frac{\partial \hat{\beta}_j}{\partial \lambda_s} \times \frac{\partial \hat{\beta}_j}{\partial \lambda_t} \end{aligned}$$

$\tilde{\ell}_i$ has second derivatives

$$\begin{aligned} \frac{\partial^2 \tilde{\ell}_i}{\partial u^2} &= \dot{\ell}_i(\tilde{u}_i(u, h)) \times \frac{\partial^2 \tilde{u}_i}{\partial u^2} + \ddot{\ell}_i(\tilde{u}_i(u, h)) \times \left(\frac{\partial \tilde{u}_i}{\partial u} \right)^2 \\ \frac{\partial^2 \tilde{\ell}_i}{\partial u \partial h} &= \dot{\ell}_i(\tilde{u}_i(u, h)) \times \frac{\partial^2 \tilde{u}_i}{\partial u \partial h} + \ddot{\ell}_i(\tilde{u}_i(u, h)) \times \frac{\partial \tilde{u}_i}{\partial u} \times \frac{\partial \tilde{u}_i}{\partial h} \\ \frac{\partial^2 \tilde{\ell}_i}{\partial h^2} &= \dot{\ell}_i(\tilde{u}_i(u, h)) \times \frac{\partial^2 \tilde{u}_i}{\partial h^2} + \ddot{\ell}_i(\tilde{u}_i(u, h)) \times \left(\frac{\partial \tilde{u}_i}{\partial h} \right)^2 \end{aligned}$$

and \tilde{u}_i has second derivatives

$$\begin{aligned}\frac{\partial^2 \tilde{u}_i}{\partial u^2} &= \frac{\ddot{\ell}_i(u) h}{(1 - \ddot{\ell}_i(u) h)^2} + \frac{(\ddot{\ell}_i(u) \ddot{\ell}_i(u) + \dot{\ell}_i(u) \ddot{\ell}_i(u)) h^2}{(1 - \ddot{\ell}_i(u) h)^2} + \frac{2 \dot{\ell}_i(u) \ddot{\ell}_i(u)^2 h^3}{(1 - \ddot{\ell}_i(u) h)^3} \\ \frac{\partial^2 \tilde{u}_i}{\partial u \partial h} &= \frac{\ddot{\ell}_i(u)}{(1 - \ddot{\ell}_i(u) h)^2} + \frac{2 \dot{\ell}_i(u) \ddot{\ell}_i(u) h}{(1 - \ddot{\ell}_i(u) h)^3} \\ \frac{\partial^2 \tilde{u}_i}{\partial h^2} &= \frac{2 \dot{\ell}_i(u) \ddot{\ell}_i(u)}{(1 - \ddot{\ell}_i(u) h)^3}\end{aligned}$$

4.1 Computational Complexity

Let q denote the number of hyperparameters. How best to proceed with the computations will depend on which is greater n or p . We assume first that $n > p$. \mathbf{H} and its Cholesky factorization \mathbf{L} can be computed in $\mathcal{O}(p^2 n)$ operations. Let \mathbf{h} denote the vector of h_i values. The complexity of computing the ALO value, gradient, and hessian is dominated by the cost of evaluating \mathbf{h} and its derivatives. Using the formula

$$h_i = \|\mathbf{L}^{-1} \mathbf{x}_i\|^2$$

\mathbf{h} can be computed with $\mathcal{O}(p^2 n)$ operations. For the derivatives of \mathbf{h} , we first compute the matrices $\frac{\partial \mathbf{H}}{\partial \lambda_s}$ which can be done in $\mathcal{O}(p^2 q n)$ operations. The derivatives can then, also, be computed with $\mathcal{O}(p^2 q n)$ operations using the formulas

$$\begin{aligned}\mathbf{t}_i &= \mathbf{L}^{-1 \top} \mathbf{L}^{-1} \mathbf{x}_i \\ \frac{\partial h_i}{\partial \lambda_s} &= \mathbf{t}_i^\top \frac{\partial \mathbf{H}}{\partial \lambda_s} \mathbf{t}_i\end{aligned}$$

For the second derivatives of \mathbf{h} , we, similarly, first compute the matrices $\frac{\partial^2 \mathbf{H}}{\partial \lambda_s \partial \lambda_t}$, which can be done in $\mathcal{O}(p^2 q^2 n)$ operations, and then compute

$$\begin{aligned}\mathbf{r}_{si} &= \mathbf{L}^{-1} \frac{\partial \mathbf{H}}{\partial \lambda_s} \mathbf{H}^{-1} \mathbf{x}_i \\ \frac{\partial^2 h_i}{\partial \lambda_s \partial \lambda_t} &= \mathbf{r}_{si}^\top \mathbf{r}_{ti} - \mathbf{t}_i^\top \frac{\partial^2 \mathbf{H}}{\partial \lambda_s \partial \lambda_t} \mathbf{t}_i\end{aligned}$$

with $\mathcal{O}(p^2 q^2 n)$ operations.

When $p > n$ and \mathbf{W} is nonsingular, we can achieve better complexity by evaluating the equations in a different order. We apply the matrix inversion lemma to $\mathbf{H}^{-1} = [\mathbf{X}^\top \mathbf{A} \mathbf{X} + \mathbf{W}]^{-1}$ to obtain

$$\mathbf{H}^{-1} = \mathbf{W}^{-1} + \mathbf{W}^{-1} \mathbf{X}^\top \left[\mathbf{A}^{-1} + \mathbf{X} \mathbf{W}^{-1} \mathbf{X}^\top \right]^{-1} \mathbf{X} \mathbf{W}^{-1}$$

Computing the matrix $\mathbf{A}^{-1} + \mathbf{X} \mathbf{W}^{-1} \mathbf{X}^\top$ and its Cholesky factorization can be done in $\mathcal{O}(n^2 p)$ operations, and a product $\mathbf{H}^{-1} \mathbf{b}$ can be done in $\mathcal{O}(np)$ operations. Applying the same approach, where we avoid explicitly computing the p -by- p matrices, we can compute the gradient and hessian in $\mathcal{O}(n^2 qp)$ and $\mathcal{O}(n^2 q^2 p)$ operations, respectively. If \mathbf{W} is singular but only has k zero diagonal entries where $k \ll p$, we can combine this approach with block matrix inversion and still achieve more efficient formulas.

5. Numerical Experiments

We run experiments designed around these lines of inquiry:

1. What do the derivatives of ALO look like?
2. How do hyperparameters found by ALO optimization compare to those found by grid search?
3. What is the cost of ALO optimization?
4. Can we use ALO optimization to fit models with multiple hyperparameters that lead to better performance on out-of-sample predictions?

To that end, we fit ridge regression and regularized logistic regression models to a range of real-world data sets. Table 1 catalogs the data sets used, and we provide brief summaries below.

Data Set	Task	n	p
Breast Cancer	classification	569	30
Cleveland Heart	classification	297	22
Pollution	regression	60	15
Arcene	classification	200	10000
Gisette	classification	7000	5000

Table 1: Sample data sets used in experiments

Breast Cancer is a binary classification data set where the objective is to predict whether breast mass is malignant from characteristics of cell nuclei.

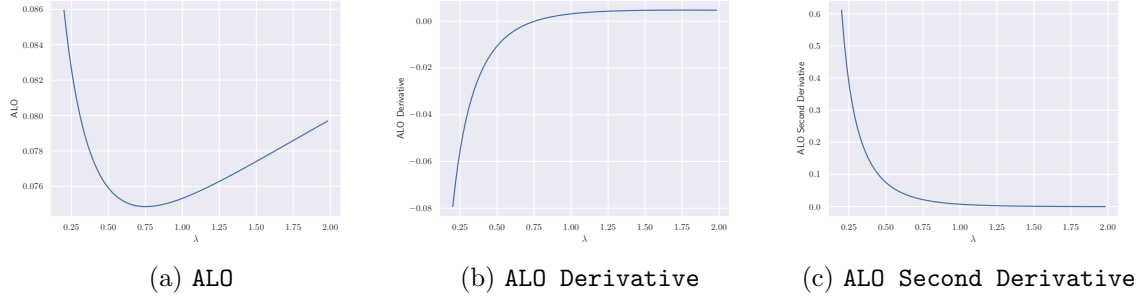
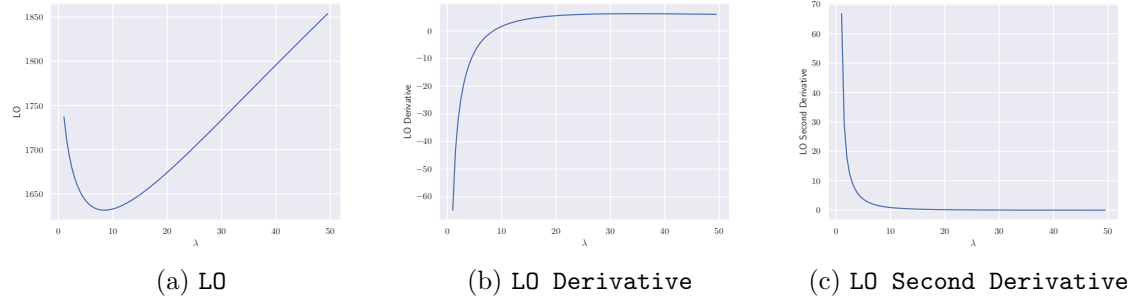
Cleveland Heart is a binary classification data set where the objective is to detect the presence of heart disease. The data set uses 13 features, but 5 are categorical. We obtain 22 features after transforming the categorical features to indicator variables, and we obtain 297 entries after dropping any entries with missing features.

Pollution is a regression data set where the objective is to predict the mortality rate of metropolitan areas from environmental and socioeconomic variables (McDonald and Schwing, 1973).

Arcene is a binary classification data set where the objective is to distinguish cancer versus normal patterns from mass-spectrometric data. We combine the **Arcene** training and validation data sets to get a data set with 200 entries.

Gisette is a binary classification data set where the objective is to predict whether a handwritten digit is a 4 or a 9. The data set includes features derived from a 28x28 pixel digit image and non-predictive probe features. For **Gisette**, we again combine the training and validation data sets to get 7000 entries.

We preprocess all data sets used for training so that features have zero mean and unit standard deviation when not constant.

Figure 1: ALO for regularized logistic regression on the **Breast Cancer** data setFigure 2: LO for ridge regression on the **Pollution** data set

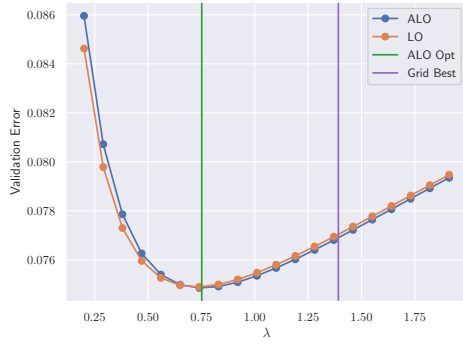
5.1 Visualizing ALO and its Derivatives

We begin by graphing ALO and its derivatives on sample data sets. Figure 1 plots ALO and its derivatives for logistic regression with the ridge regularizer $R_\lambda(\beta) = \lambda \|\beta\|^2$ on the **Breast Cancer** data set. Figure 2 plots LO and its derivatives for ridge regression on the **Pollution** data set. With ridge regression, ALO and LO are equivalent.

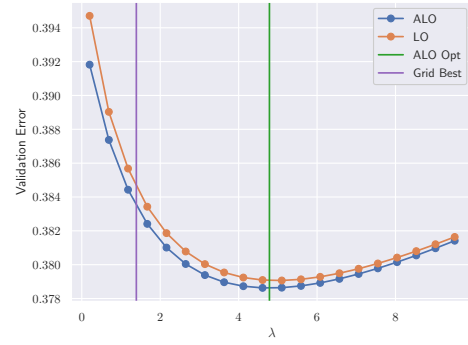
5.2 Comparing ALO Optimization to Grid Search

We next compare hyperparameters found by ALO optimization to those found by grid search for the ridge regularizer. For ALO optimization, we use the Python package **peak-engines** available from <https://github.com/rnburn/peak-engines>; and for grid search, we use **LogisticRegressionCV** and **RidgeCV** from **sklearn-0.23.1** with default settings. **LogisticRegressionCV** defaults to run a grid search with a 5-fold cross-validation and 10 points of evaluation. **RidgeCV** defaults to run a grid search with LO and 3 points of evaluation.¹ Figure 3 shows where the hyperparameters found by grid search and ALO optimization lie along the LO and ALO curves. We additionally benchmark how long it takes ALO optimization and grid search to find hyperparameters on the sample data sets. Figure 4 shows the durations measured in seconds.

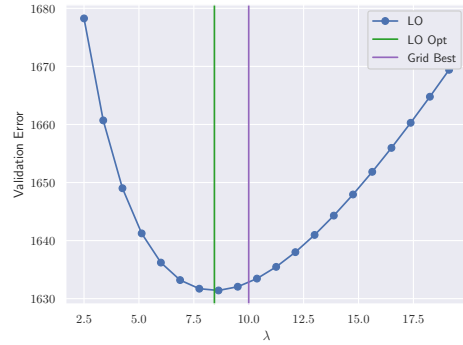
1. **RidgeCV**'s documentation claims that it uses Generalized Cross-Validation, but it's really using LO. See <https://github.com/scikit-learn/scikit-learn/issues/18079>.



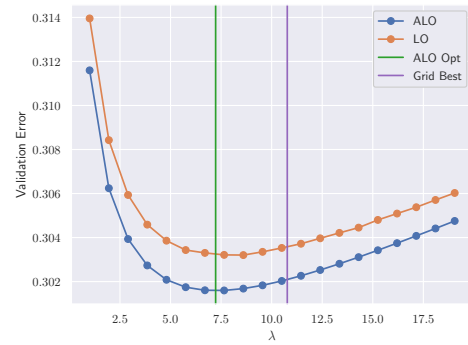
(a) Breast Cancer



(b) Cleveland Heart



(c) Pollution



(d) Arcene

Figure 3: Plot of ALO and LO curves for sample data sets. The green line shows the hyperparameter found by ALO optimization, and the purple line shows the hyperparameter found by grid search.

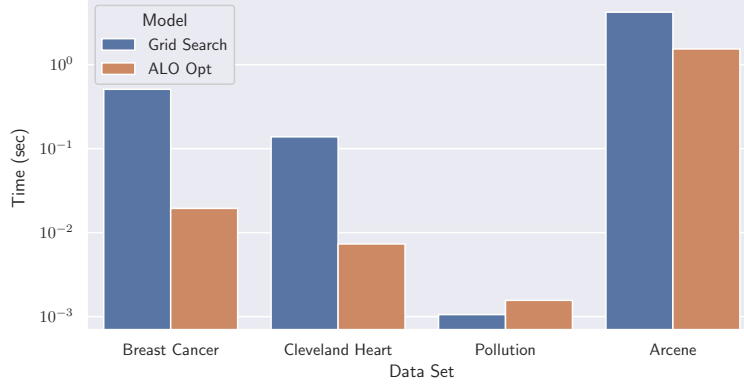


Figure 4: Benchmark showing how long ALO optimization and grid search take to find hyperparameters on sample data sets. The times are averaged over 10 runs and shown on a logarithmic scale.

5.3 Fitting Models with Multiple Hyperparameters

To test ALO optimization with multiple hyperparameters, we fit logistic regression with a bridge regularizer (Fu , 1998) of the form $\lambda_1 |\beta_j|^{\lambda_2}$ where $\lambda_2 \geq 1$ to a 5-fold cross-validation of the **Gisette** data set. To ensure that $r_{j,\lambda}$ has a continuous fourth derivative, we interpolate the regularizer with a polynomial when $|\beta_j| < \delta$ so that

$$r_{j,\lambda}(\beta_j) = \begin{cases} \lambda_1 |\beta_j|^{\lambda_2} & \text{if } |\beta_j| \geq \delta \\ \lambda_1 p_{\lambda_2}(|\beta_j|) & \text{if } |\beta_j| < \delta \end{cases}$$

where $p_{\lambda_2}(t) = a_1 t^2 + a_2 t^4 + a_3 t^5 + a_4 t^6 + a_5 t^7$ and \mathbf{a} is chosen so that $p_{\lambda_2}(\delta)$ matches δ^{λ_2} up to the fourth derivative. For our experiment, we use $\delta = 0.01$. For comparison, we also fit logistic regression with ridge regularization. Table 2 shows the hyperparameters found for bridge regularization and ridge regularization and compares their performance on each cross-validation fold.

6. Discussion

In this paper, we demonstrated how to select hyperparameters by computing the gradient and hessian of ALO and applying a second-order optimizer to find a local minimum. The approach is applicable to a large class of commonly used models, including regularized logistic regression and ridge regression. We applied ALO optimization to fit regularized models to a variety of different real-world data sets. We found that when using a single-parameter regularizer, we were able to find hyperparameters with better LO scores than standard grid-search approaches and frequently were able to do so in less time. ALO optimization, furthermore, scales to handle multiple hyperparameters, and we demonstrated how it could be used to fit hyperparameters for bridge regularization.

Fold	λ_1	λ_2	Test Error	Fold	λ	Test Error
0	17.3	1.95	0.0698	0	19.1	0.0709
1	16.8	1.94	0.0581	1	18.7	0.0587
2	13.0	1.81	0.0667	2	18.3	0.0697
3	10.5	1.75	0.0677	3	16.5	0.0714
4	10.8	1.83	0.0639	4	14.4	0.0667
Mean			0.0652	Mean		0.0675

(a) Bridge Regularization

(b) Ridge Regularization

Table 2: Hyperparameters and performance of ALO optimization with logistic regression using bridge and ridge regularization on a 5-fold cross-validation of the **Gisette** data set. Test error is measured as the negative mean log-likelihood of the out-of-sample fold.

Acknowledgments

We made use of the UCI Machine Learning Repository <http://archive.ics.uci.edu/ml> in our experiments.

Appendix A: Proof of Theorem 1

To derive the derivative of $\hat{\beta}$, we observe that at the optimum of Equation 1 the gradient is zero.

$$\sum_{i=1}^n \dot{\ell}_i(u_i) \mathbf{x}_i + \nabla R_{\lambda}(\hat{\beta}) = 0$$

Differentiating both sides of the equation gives us

$$\begin{aligned} \frac{\partial}{\partial \lambda_s} \left(\sum_{i=1}^n \dot{\ell}_i(u_i) \mathbf{x}_i + \nabla R_{\lambda}(\hat{\beta}) \right) &= 0 \\ \Leftrightarrow \mathbf{H} \left(\frac{\partial \hat{\beta}}{\partial \lambda_s} \right) + \frac{\partial \nabla R_{\lambda}}{\partial \lambda_s}(\hat{\beta}) &= 0 \\ \Leftrightarrow \frac{\partial \hat{\beta}}{\partial \lambda_s} &= -\mathbf{H}^{-1} \frac{\partial \nabla R_{\lambda}}{\partial \lambda_s}(\hat{\beta}) \end{aligned}$$

For the derivative of h_i , we apply this formula for differentiating an inverse matrix

$$\frac{d\mathbf{E}^{-1}}{dt} = -\mathbf{E}^{-1} \frac{d\mathbf{E}}{dt} \mathbf{E}^{-1}$$

The other derivatives are derived as straightforward differentiations of their associated value formulas.

Appendix B: Proof of Theorem 2

For the second derivative of $\hat{\beta}$, we derive

$$\begin{aligned}\frac{\partial^2 \hat{\beta}}{\partial \lambda_s \partial \lambda_t} &= \frac{\partial}{\partial \lambda_t} \left[-\mathbf{H}^{-1} \frac{\partial \nabla R_{\lambda}}{\partial \lambda_s}(\hat{\beta}) \right] \\ &= \mathbf{H}^{-1} \frac{\partial \mathbf{H}}{\partial \lambda_t} \mathbf{H}^{-1} \frac{\partial \nabla R_{\lambda}}{\partial \lambda_s}(\hat{\beta}) - \mathbf{H}^{-1} \frac{\partial}{\partial \lambda_t} \left(\frac{\partial \nabla R_{\lambda}}{\partial \lambda_s}(\hat{\beta}) \right) \\ &= -\mathbf{H}^{-1} \frac{\partial \mathbf{H}}{\partial \lambda_t} \frac{\partial \hat{\beta}}{\partial \lambda_s} - \mathbf{H}^{-1} \frac{\partial}{\partial \lambda_t} \left(\frac{\partial \nabla R_{\lambda}}{\partial \lambda_s}(\hat{\beta}) \right)\end{aligned}$$

Now,

$$\begin{aligned}\left(\frac{\partial \mathbf{H}}{\partial \lambda_t} \right) \frac{\partial \hat{\beta}}{\partial \lambda_s} &= \left[\mathbf{X}^{\top} \frac{\partial \mathbf{A}}{\partial \lambda_t} \mathbf{X} + \frac{\partial \mathbf{W}}{\partial \lambda_t} \right] \frac{\partial \hat{\beta}}{\partial \lambda_s} \\ &= \mathbf{X}^{\top} \cdot \text{vec} \left[\left\{ \ddot{\ell}_i(u_i) \times \frac{\partial u_i}{\partial \lambda_s} \times \frac{\partial u_i}{\partial \lambda_t} \right\}_i \right] + \frac{\partial \nabla^2 R_{\lambda}}{\partial \lambda_t}(\hat{\beta}) \frac{\partial \hat{\beta}}{\partial \lambda_s} \\ &\quad + \text{vec} \left[\left\{ \ddot{r}_{j,\lambda}(\hat{\beta}_j) \times \frac{\partial \hat{\beta}_j}{\partial \lambda_s} \times \frac{\partial \hat{\beta}_j}{\partial \lambda_t} \right\}_j \right]\end{aligned}$$

and

$$\frac{\partial}{\partial \lambda_t} \left(\frac{\partial \nabla R_{\lambda}}{\partial \lambda_s}(\hat{\beta}) \right) = \frac{\partial \nabla^2 R_{\lambda}}{\partial \lambda_s}(\hat{\beta}) \frac{\partial \hat{\beta}}{\partial \lambda_t} + \frac{\partial \nabla R_{\lambda}}{\partial \lambda_s \partial \lambda_t}(\hat{\beta})$$

Combining the equations, the second derivative of $\hat{\beta}$ becomes

$$\begin{aligned}\frac{\partial^2 \hat{\beta}}{\partial \lambda_s \partial \lambda_t} &= -\mathbf{H}^{-1} \mathbf{X}^{\top} \cdot \text{vec} \left[\left\{ \ddot{\ell}_i(u_i) \times \frac{\partial u_i}{\partial \lambda_s} \times \frac{\partial u_i}{\partial \lambda_t} \right\}_i \right] \\ &\quad - \mathbf{H}^{-1} \frac{\partial \nabla^2 R_{\lambda}}{\partial \lambda_s}(\hat{\beta}) \frac{\partial \hat{\beta}}{\partial \lambda_t} - \mathbf{H}^{-1} \frac{\partial \nabla^2 R_{\lambda}}{\partial \lambda_t}(\hat{\beta}) \frac{\partial \hat{\beta}}{\partial \lambda_s} - \mathbf{H}^{-1} \frac{\partial \nabla R_{\lambda}}{\partial \lambda_s \partial \lambda_t}(\hat{\beta}) \\ &\quad - \mathbf{H}^{-1} \cdot \text{vec} \left[\left\{ \ddot{r}_{j,\lambda}(\hat{\beta}_j) \times \frac{\partial \hat{\beta}_j}{\partial \lambda_s} \times \frac{\partial \hat{\beta}_j}{\partial \lambda_t} \right\}_j \right]\end{aligned}$$

For the second derivative of h_i , we derive

$$\begin{aligned}\frac{\partial^2 h_i}{\partial \lambda_s \partial \lambda_t} &= \frac{\partial}{\partial \lambda_t} \left(-\mathbf{x}_i^{\top} \mathbf{H}^{-1} \frac{\partial \mathbf{H}}{\partial \lambda_s} \mathbf{H}^{-1} \mathbf{x}_i \right) \\ &= \mathbf{x}_i^{\top} \mathbf{H}^{-1} \frac{\partial \mathbf{H}}{\partial \lambda_t} \mathbf{H}^{-1} \frac{\partial \mathbf{H}}{\partial \lambda_s} \mathbf{H}^{-1} \mathbf{x}_i + \mathbf{x}_i^{\top} \mathbf{H}^{-1} \frac{\partial \mathbf{H}}{\partial \lambda_s} \mathbf{H}^{-1} \frac{\partial \mathbf{H}}{\partial \lambda_t} \mathbf{H}^{-1} \mathbf{x}_i \\ &\quad - \mathbf{x}_i^{\top} \mathbf{H}^{-1} \frac{\partial^2 \mathbf{H}}{\partial \lambda_s \partial \lambda_t} \mathbf{H}^{-1} \mathbf{x}_i \\ &= 2\mathbf{x}_i^{\top} \mathbf{H}^{-1} \frac{\partial \mathbf{H}}{\partial \lambda_s} \mathbf{H}^{-1} \frac{\partial \mathbf{H}}{\partial \lambda_t} \mathbf{H}^{-1} \mathbf{x}_i - \mathbf{x}_i^{\top} \mathbf{H}^{-1} \frac{\partial^2 \mathbf{H}}{\partial \lambda_s \partial \lambda_t} \mathbf{H}^{-1} \mathbf{x}_i\end{aligned}$$

and

$$\begin{aligned}
\frac{\partial^2 \mathbf{H}}{\partial \lambda_s \partial \lambda_t} &= \frac{\partial}{\partial \lambda_t} \left(\mathbf{X}^\top \frac{\partial \mathbf{A}}{\partial \lambda_s} \mathbf{X} + \frac{\partial \mathbf{W}}{\partial \lambda_s} \right) \\
&= \mathbf{X}^\top \frac{\partial^2 \mathbf{A}}{\partial \lambda_s \partial \lambda_t} \mathbf{X} + \frac{\partial^2 \mathbf{W}}{\partial \lambda_s \partial \lambda_t} \\
\left(\frac{\partial^2 \mathbf{W}}{\partial \lambda_s \partial \lambda_t} \right)_{jj} &= \frac{\partial}{\partial \lambda_t} \left(\frac{\partial \ddot{r}_{j,\lambda}}{\partial \lambda_s}(\hat{\beta}_j) + \ddot{r}_{j,\lambda}(\hat{\beta}_j) \times \frac{\partial \hat{\beta}_j}{\partial \lambda_s} \right) \\
&= \frac{\partial^2 \ddot{r}_{j,\lambda}}{\partial \lambda_s \partial \lambda_t}(\hat{\beta}_j) + \frac{\partial \ddot{r}_{j,\lambda}}{\partial \lambda_s}(\hat{\beta}_j) \times \frac{\partial \hat{\beta}_j}{\partial \lambda_t} + \frac{\partial \ddot{r}_{j,\lambda}}{\partial \lambda_t}(\hat{\beta}_j) \times \frac{\partial \hat{\beta}_j}{\partial \lambda_s} \\
&\quad + \ddot{r}_{j,\lambda}(\hat{\beta}_j) \times \frac{\partial^2 \hat{\beta}_j}{\partial \lambda_s \partial \lambda_t} + \ddot{r}_{j,\lambda}(\hat{\beta}_j) \times \frac{\partial \hat{\beta}_j}{\partial \lambda_s} \times \frac{\partial \hat{\beta}_j}{\partial \lambda_t}
\end{aligned}$$

We omit the steps for the other derivatives as they are straightforward.

Appendix C: Special Cases

We provide formulas for the special cases of ridge regression and logistic regression with ridge regularization. We parameterize ridge regularization so that

$$R_\lambda(\boldsymbol{\beta}) \triangleq \lambda^2 \|\boldsymbol{\beta}\|^2$$

With this parameterization, R_λ is defined for all λ .

Corollary 3 *For logistic regression with ridge regularization, terms in Theorem 1 simplify to*

$$\begin{aligned}
\frac{\partial \hat{\boldsymbol{\beta}}}{\partial \lambda} &= -4\lambda \mathbf{H}^{-1} \hat{\boldsymbol{\beta}} \\
\frac{\partial \mathbf{H}}{\partial \lambda} &= \mathbf{X}^\top \frac{\partial \mathbf{A}}{\partial \lambda} \mathbf{X} + 4\lambda \mathbf{I}
\end{aligned}$$

with

$$\begin{aligned}
\ell_i(u) &\triangleq \log[1 + \exp(-y_i u)] \\
\dot{\ell}_i(u) &= \frac{-y_i}{1 + \exp(y_i u)} \\
\ddot{\ell}_i(u) &= qp \\
\ddot{\ell}_i(u) &= qp(q - p)
\end{aligned}$$

where $p = (1 + \exp(-u))^{-1}$ and $q = 1 - p$.

Corollary 4 *For logistic regression with ridge regularization, terms in Theorem 2 simplify to*

$$\begin{aligned}\frac{\partial^2 \hat{\boldsymbol{\beta}}}{\partial \lambda^2} &= -\mathbf{H}^{-1} \mathbf{X}^\top \cdot \text{vec} \left[\left\{ \ddot{\ell}_i(u_i) \times \left(\frac{\partial u_i}{\partial \lambda} \right)^2 \right\}_i \right] \\ &\quad - 8\lambda \mathbf{H}^{-1} \frac{\partial \hat{\boldsymbol{\beta}}}{\partial \lambda} - 4\mathbf{H}^{-1} \hat{\boldsymbol{\beta}} \\ \frac{\partial^2 \mathbf{H}}{\partial \lambda^2} &= \mathbf{X}^\top \frac{\partial^2 \mathbf{A}}{\partial \lambda^2} \mathbf{X} + 4\mathbf{I}\end{aligned}$$

with

$$\ddot{\ell}_i(u) = qp(q^2 + p^2) - 4q^2p^2$$

where $p = (1 + \exp(-u))^{-1}$ and $q = 1 - p$.

Corollary 5 *For ridge regression, terms in Theorem 1 simplify to*

$$\begin{aligned}\frac{\partial \hat{\boldsymbol{\beta}}}{\partial \lambda} &= -4\lambda \mathbf{H}^{-1} \hat{\boldsymbol{\beta}} \\ \frac{\partial h_i}{\partial \lambda} &= -4\lambda \mathbf{x}_i^\top \mathbf{H}^{-1} \mathbf{H}^{-1} \mathbf{x}_i\end{aligned}$$

and

$$\begin{aligned}\frac{\partial \tilde{u}_i}{\partial u} &= \frac{1}{1 - 2h} \\ \frac{\partial \tilde{u}_i}{\partial h} &= \frac{\dot{\ell}_i(u)}{(1 - 2h)^2}\end{aligned}$$

with

$$\begin{aligned}\ell_i(u) &= (u - y_i)^2 \\ \dot{\ell}_i(u) &= 2(u - y_i)\end{aligned}$$

Corollary 6 *For ridge regression, terms in Theorem 2 simplify to*

$$\begin{aligned}\frac{\partial^2 \hat{\boldsymbol{\beta}}}{\partial \lambda^2} &= -8\lambda \mathbf{H}^{-1} \frac{\partial \hat{\boldsymbol{\beta}}}{\partial \lambda} - 4\mathbf{H}^{-1} \hat{\boldsymbol{\beta}} \\ \frac{\partial^2 h_i}{\partial \lambda^2} &= 32\lambda^2 \mathbf{x}_i^\top \mathbf{H}^{-1} \mathbf{H}^{-1} \mathbf{H}^{-1} \mathbf{x}_i - 4\mathbf{x}_i^\top \mathbf{H}^{-1} \mathbf{H}^{-1} \mathbf{x}_i\end{aligned}$$

and

$$\begin{aligned}
\frac{\partial^2 \tilde{\ell}_i}{\partial u^2} &= 2 \left(\frac{\partial \tilde{u}_i}{\partial u} \right)^2 \\
\frac{\partial^2 \tilde{\ell}_i}{\partial u \partial h} &= \dot{\ell}_i(\tilde{u}_i(u, h)) \times \frac{\partial^2 \tilde{u}_i}{\partial u \partial h} + 2 \times \frac{\partial \tilde{u}_i}{\partial u} \times \frac{\partial \tilde{u}_i}{\partial h} \\
\frac{\partial^2 \tilde{\ell}_i}{\partial h^2} &= \dot{\ell}_i(\tilde{u}_i(u, h)) \times \frac{\partial^2 \tilde{u}_i}{\partial h^2} + 2 \left(\frac{\partial \tilde{u}_i}{\partial h} \right)^2 \\
\frac{\partial^2 \tilde{u}_i}{\partial u \partial h} &= \frac{2}{(1-2h)^2} \\
\frac{\partial^2 \tilde{u}_i}{\partial h^2} &= \frac{4\dot{\ell}_i(u)}{(1-2h)^3}
\end{aligned}$$

References

- D. M. Allen. The relationship between variable selection and data augmentation and a method for prediction. *Technometrics*, 16:125–127, 1974.
- Y. Bengio. Gradient-based optimization of hyperparameters. *Neural Computation*, 12(8):1889–1900, 2000.
- J. Bergstra and Y. Bengio. Random search for hyper-parameter optimization. *Journal of Machine Learning Research*, 13:281–305, 2012.
- C. B. Do, D. A. Woods, A. Y. Ng. Efficient multiple hyperparameter learning for log-linear models. *Advances in Neural Information Processing Systems*, 20:377–384, 2008.
- W. J. Fu. Penalized regression: the bridge versus the lasso. *Journal of Computational and Graphical Statistics*, 7:397–416, 1998.
- G. C. McDonald and R. C. Schwing. Instabilities of regression estimates relating air pollution to mortality. *Technometrics*, 15:463–481, 1973.
- J. J. Moré and D. C. Sorensen. Computing a trust region step. *SIAM Journal on Scientific and Statistical Computing*, 4(3):553–572, 1983.
- J. Nocedal and S. J. Wright. Numerical Optimization. *Springer Series in Operations Research and Financial Engineering*. Springer, Second edition, 1999.
- K. R. Rad and A. Maleki. A scalable estimate of the extra-sample prediction error via approximate leave-one-out. *ArXiv:1801.10243v4*, 2020.
- K. R. Rad, W. Zhou, A. Maleki. Error bounds in estimating the out-of-sample prediction error using leave-one-out cross validation in high-dimensions. *ArXiv:2003.01770v1*, 2020.
- D. C. Sorensen. Newton’s method with a model trust region modification. *SIAM Journal on Numerical Analysis*, 19(2):409–426, 1982.