

3 Regularized Linear Regression Using Lasso [14 Points]

Lasso is a form of regularized linear regression, where the L1 norm of the parameter vector is penalized. It is used in an attempt to get a sparse parameter vector where features of little “importance” are assigned to zero weight. But why does lasso encourage sparse parameters? For this question, you are going to examine this.

Let \mathbf{X} denote an $n \times d$ matrix where rows are training points, \mathbf{y} denotes an $n \times 1$ vector of corresponding output value, \mathbf{w} denotes a $d \times 1$ parameter vector and \mathbf{w}^* denotes the optimal parameter vector. To make the analysis easier we will consider the special case where the training data is whitened (i.e., $\mathbf{X}^\top \mathbf{X} = \mathbf{I}$). For lasso regression, the optimal parameter vector is given by

$$\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmin}} J_\lambda(\mathbf{w}), \quad (1)$$

where $J_\lambda(\mathbf{w})$ is the function we want to minimize, which is given by

$$J_\lambda(\mathbf{w}) = \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_1, \quad (2)$$

where $\lambda > 0$. Note that the L_1 norm for a vector $\mathbf{w} = [w_1, \dots, w_d]^\top \in \mathbb{R}^d$ is defined as $\|\mathbf{w}\|_1 = |w_1| + \dots + |w_d|$.

1. [3 Points] In 3.2 and 3.3, we will show that whitening the training data nicely decouples the features, making \mathbf{w}_i^* determined by the i th feature and the output regardless of other features. To show this, begin by writing $J_\lambda(\mathbf{w})$ in the form

$$J_\lambda(\mathbf{w}) = g(\mathbf{y}) + \sum_{i=1}^d f(X_{\cdot i}, \mathbf{y}, w_i, \lambda), \quad (3)$$

where $X_{\cdot i}$ is the i th column of \mathbf{X} , g is a function of only \mathbf{y} and f is a function of $X_{\cdot i}, \mathbf{y}, w_i, \lambda$

$$\begin{aligned} J_\lambda(\mathbf{w}) &= \frac{1}{2} (\vec{\mathbf{y}} - \vec{\mathbf{X}} \vec{\mathbf{w}})^\top (\vec{\mathbf{y}} - \vec{\mathbf{X}} \vec{\mathbf{w}}) + \lambda \sum |w_i| \\ &= \frac{1}{2} [\vec{\mathbf{y}}^\top \vec{\mathbf{y}} - \vec{\mathbf{y}}^\top (\vec{\mathbf{X}} \vec{\mathbf{w}}) - (\vec{\mathbf{X}} \vec{\mathbf{w}})^\top \vec{\mathbf{y}} + (\vec{\mathbf{X}} \vec{\mathbf{w}})^\top (\vec{\mathbf{X}} \vec{\mathbf{w}})] + \lambda \sum |w_i| \\ &= \frac{1}{2} (\vec{\mathbf{y}}^\top \vec{\mathbf{y}}) + \frac{1}{2} (\vec{\mathbf{X}} \vec{\mathbf{w}})^\top (\vec{\mathbf{X}} \vec{\mathbf{w}}) - \vec{\mathbf{y}}^\top (\vec{\mathbf{X}} \vec{\mathbf{w}}) + \lambda \sum |w_i| \\ &= \underbrace{\frac{1}{2} (\vec{\mathbf{y}}^\top \vec{\mathbf{y}})}_{g(\vec{\mathbf{y}})} + \sum_{i=1}^d \underbrace{\left[\frac{1}{2} (X_{\cdot i} \vec{\mathbf{w}})^\top (X_{\cdot i} \vec{\mathbf{w}}) - \vec{\mathbf{y}}^\top (X_{\cdot i} \vec{\mathbf{w}}) + \lambda |w_i| \right]}_{f(X_{\cdot i}, \vec{\mathbf{y}}, w_i, \lambda)} \quad \star \end{aligned}$$

- p is the dimension of data (number of features in each observation)

$$\log \frac{P(y_i=1 | x_i, w)}{1 - P(y_i=1 | x_i, w)} = w^T x_i$$

$$\frac{P(y_i=1 | x_i, w)}{1 - P(y_i=1 | x_i, w)} = e^{w^T x_i} = \frac{e^{w^T x_i}}{1 + e^{w^T x_i}} \cdot \frac{1}{\frac{1}{1 + e^{w^T x_i}}}$$

$$L(w) = \log P(y | X, w)$$

$$P(y | X, w) = \prod_{i=1}^n (P(y_i=1 | x_i, w))^{y_i} (1 - P(y_i=1 | x_i, w))^{(1-y_i)}$$

$$\Rightarrow \log P(y | X, w) = \sum_{i=1}^n [y_i \log(P(y_i=1 | x_i, w)) + (1-y_i) \log(1 - P(y_i=1 | x_i, w))]$$

$$= \sum_{i=1}^n [y_i \log \frac{P(y_i=1 | x_i, w)}{1 - P(y_i=1 | x_i, w)} + \log(1 - P(y_i=1 | x_i, w))]$$

$$= \sum_{i=1}^n [y_i \overbrace{\log \frac{1}{1 + e^{-w^T x_i}}}^{\log 1 - \log(1 + e^{-w^T x_i})} + \log(1 - P(y_i=1 | x_i, w))]$$

$$= \sum_{i=1}^n [y_i w^T x_i - \log(1 + e^{w^T x_i})]$$

4.2 Convex Optimization

Our goal is to find the weight vector w that maximizes this likelihood. Unfortunately, for this model, we cannot derive a closed-form solution with MLE. An alternative way to solve for w is to use gradient ascent, and update w step by step towards the optimal w . But we know gradient ascent will converge to the optimal solution w that maximizes the conditional log likelihood \mathcal{L} when \mathcal{L} is concave. In this question, you will prove that \mathcal{L} is indeed a concave function.

1. [3 points] A real-valued function $f : S \rightarrow \mathcal{R}$ defined on a convex set S , is said to be *convex* if,

$$f(tx_1 + (1-t)x_2) \leq tf(x_1) + (1-t)f(x_2), \forall x_1, x_2 \in S, \forall t \in [0, 1].$$

Show that a linear combination of n convex functions, f_1, f_2, \dots, f_n , $\sum_{i=1}^n a_i f_i(x)$ is also a convex function $\forall a_i \in \mathcal{R}^+$.

$$\begin{aligned} \sum_{i=1}^n a_i f_i(x) &= \sum_{i=1}^n a_i f_i(tx_1 + (1-t)x_2) \\ &= a_1 f_1(tx_1 + (1-t)x_2) + \dots + a_n f_n(tx_1 + (1-t)x_2) \\ \underline{a_i \in \mathcal{R}^+} &\rightarrow \leq a_1 (tf_1(x_1) + (1-t)f_1(x_2)) + \dots + \\ &\quad a_n (tf_n(x_1) + (1-t)f_n(x_2)) \\ &= t(a_1 f_1(x_1) + \dots + a_n f_n(x_1)) \\ &\quad + (1-t)(a_1 f_1(x_2) + \dots + a_n f_n(x_2)) \\ &= \sum_{i=1}^n a_i (tf_i(x_1) + (1-t)f_i(x_2)) \\ f_1, f_2, \dots, f_n \text{ are convex} &\Rightarrow \sum_{i=1}^n f_i \text{ is convex} \\ a_i \text{ is nonnegative} & \\ \sum_{i=1}^n a_i f_i(x) &\text{ is also convex.} \end{aligned}$$

2. [2 points] Show that a linear combination of n concave functions, f_1, f_2, \dots, f_n , $\sum_{i=1}^n a_i f_i(x)$ is also a concave function $\forall a_i \in \mathbb{R}^+$. Recall that if a function $f(x)$ is convex, then $-f(x)$ is concave. (You can use the result from part (1))

$f(x)$ is convex when

$$f(tx_1 + (1-t)x_2) \leq tf(x_1) + (1-t)f(x_2)$$

$$\xrightarrow{x(t)} -f(tx_1 + (1-t)x_2) \geq -tf(x_1) - (1-t)f(x_2)$$

Let $g(x) = -f(x)$ is concave

$$\Rightarrow g(tx_1 + (1-t)x_2) \geq tg(x_1) + (1-t)g(x_2)$$

$\Rightarrow g(x)$ is concave

\Rightarrow the sum of n concave function
with a nonnegative a_i ($\sum_{i=1}^n a_i f_i(x)$)

is also concave.

3. [4 points] Another property of twice differentiable convex functions is that the second derivative is non-negative. Using this property, show that $f(x) = \log(1 + \exp x)$ is a convex function. Note that this property is both sufficient and necessary. i.e. (if $f''(x)$ exists, then $f''(x) \geq 0 \iff f$ is convex)

$$f(x) = \log(1 + e^x)$$

$$\xrightarrow{d} f'(x) = \frac{1}{1+e^x} \times e^x = e^x(1+e^x)^{-1}$$

$$\begin{aligned} \xrightarrow{d} f''(x) &= e^x(1+e^x)^{-1} - (e^x)^2(1+e^x)^{-2} \\ &= \frac{e^{2x} + e^x - e^{2x}}{(1+e^x)^2} = \frac{e^x}{(1+e^x)^2} \text{ (exists)} \end{aligned}$$

$$\Rightarrow f''(x) = \frac{e^x}{(1+e^x)^2} \geq 0$$

$$\Rightarrow f \text{ is convex} \quad \#$$

4. [4 points] Let $f_i : \mathcal{S} \rightarrow \mathcal{R}$ for $i = 1, \dots, n$ be a set of convex functions. Is $f(x) = \max_i f_i(x)$ also convex? If yes, prove it. If not, provide a counterexample.

Yes

$$f(t x_1 + (1-t)x_2) = f_i(t x_1 + (1-t)x_2) \text{ for } i = 1, \dots, n$$

$\because f_i$ is convex

$$\Rightarrow f_i(t x_1 + (1-t)x_2) \leq t f_i(x_1) + (1-t)f_i(x_2)$$

$$\leq t \max_i f_i(x_1) + (1-t) \max_i f_i(x_2)$$

$$= t f(x_1) + (1-t)f(x_2)$$

$$\Rightarrow f(t x_1 + (1-t)x_2) \leq t f(x_1) + (1-t)f(x_2)$$

$$\Rightarrow f(x) = \max_i f_i(x) \text{ also convex}$$