

LASSO Notes

Justin - Ds



LASSO

LASSO (Least Absolute Shrinkage and Selection Operator) estimates β^* by solving the following convex optimization problem:

$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^d} \|Y - X\beta\|_2^2 + \lambda \|\beta\|_1,$$

where:

- $\|Y - X\beta\|_2^2$: residual sum of squares (RSS). 管差平方和
- $\|\beta\|_1 = \sum_{j=1}^d |\beta_j|$: ℓ_1 -norm penalty.
- $\lambda > 0$: tuning parameter that controls the trade-off between goodness of fit and sparsity. 拟合“好坏”和“稀疏性” balance

Not
proof

- Interpretation as MAP estimator with a Laplace prior on β^*

- Questions:

- How to compute LASSO estimate?
- What is the statistical properties of LASSO?

chatgpt

LASSO 亂世斯爾

How to compute LASSO: proximal gradient method

LASSO 可看作 MAP of prior func of $p(\beta_i^*) = \frac{\tau}{2} e^{-\tau|\beta_i^*|}$

β_i^* ~ density of Laplace distribution

$$\tau = \frac{\lambda}{2\sigma}$$

A more general class of convex optimization

Consider unconstrained convex optimization problem of the form 凸优化问题

$$\min_{\mathbf{x} \in \mathbb{R}^d} F(\mathbf{x}) := f(\mathbf{x}) + h(\mathbf{x})$$

where

- $f(\mathbf{x})$: a differentiable, convex function 可微凸函数
- $h(\mathbf{x})$: a convex, potentially non-differentiable function (e.g., ℓ_1 -norm). 不确定
- Example: LASSO can be viewed as taking 可微性

$$f(\mathbf{x}) = \|\mathbf{Y} - \mathbf{X}\beta\|_2^2, \quad h(\mathbf{x}) = \lambda \|\beta\|_1.$$

可微 凸但不可微

Issue: gradient descent (GD) does not work (due to non-smoothness)

在 $\beta_j=0$ 时不可微，不能用梯度下降

处理非光滑函数，对 $h(x)$ 通过近端映射 (proximal mapping) A Proximal View of Gradient Descent

近端梯度法

- To motivate proximal gradient methods, we first revisit gradient descent for $\min_x f(x)$, where $f(\cdot)$ is convex and smooth
- Gradient descent update: $x_{t+1} = x_t - \eta \nabla f(x_t)$ 对 f 以梯度下降
- This is equivalent to

$$x_{t+1} = \arg \min_x \left\{ \underbrace{f(x_t) + \langle \nabla f(x_t), x - x_t \rangle}_{\text{first-order approximation at } x_t} + \underbrace{\frac{1}{2\eta} \|x - x_t\|_2^2}_{\text{proximal term}} \right\}$$

近端项

- Heuristics: search for x_{t+1} that 在 x_t 处一阶泰勒展开线性近似
 - aim to minimize $f(\cdot)$ (through minimizing first-order approximation)
 - remains close to x_t such that first-order approximation at x_t is valid (enforced by proximal term)
- Benefit: minimizing a quadratic function, admits simple solution (i.e., GD)
加入对步长的限制，且最小化 $f(\cdot)$

Proximal gradient method: algorithm

Consider an iterative algorithm: starting from x_t , update

$$x_{t+1} = \arg \min_x \left\{ f(x_t) + \underbrace{\langle \nabla f(x_t), x - x_t \rangle}_{\text{first-order approximation at } x_t} + h(x) + \underbrace{\frac{1}{2\eta} \|x - x_t\|_2^2}_{\text{proximal term}} \right\}$$

平衡更新方向
加速收敛 / 保证凸优化
收敛性

不可微凸
函数
(处理正则化
也简单)

限制更新步长
近端梯度

- Define proximal operator

近端算子

$$\text{prox}_h(v) = \arg \min_{x \in \mathbb{R}^d} \left\{ h(x) + \frac{1}{2} \|x - v\|_2^2 \right\}$$

- If this proximal operator is easy to compute, then we can express

$$x_{t+1} = \text{prox}_{\eta h}(x_t - \eta \nabla f(x_t))$$

- alternates between gradient updates on f and proximal minimization on h

首先梯度下降

然后通过近端算子调整

Proximal view of (GD): (梯度下降法)

Proof: $\ell(x) = f(x_t) + \langle \nabla f(x_t), x - x_t \rangle + \frac{1}{2\eta} \|x - x_t\|_2^2$

$$= \frac{1}{2\eta} x^T x - \frac{1}{\eta} \langle x, x_t \rangle + \langle \nabla f(x_t), x \rangle + \text{const.}$$

$$\nabla \ell(x) = \frac{1}{\eta} x - \frac{1}{\eta} x_t + \nabla f(x_t)$$

$$= \frac{1}{\eta} (x - x_t + \eta \cdot \nabla f(x_t)) = 0$$

$$\Rightarrow x = x_t - \eta \nabla f(x_t)$$

x here is minimizer of $\ell(x)$

Proximal gradient method [正端梯度法]

$$\begin{aligned}\ell(x) &= f(x_t) + \langle \nabla f(x_t), x - x_t \rangle + h(x) + \frac{1}{2\eta} \|x - x_t\|_2^2 \\ &= h(x) + \langle \nabla f(x_t), x - x_t \rangle + \frac{1}{2\eta} \|x - x_t\|_2^2 + \text{const.} \\ &= h(x) + \frac{1}{2\eta} \|x - x_t + \eta \nabla f(x_t)\|_2^2 - \frac{\eta}{2} \|\nabla f(x_t)\|_2^2 + \text{const.} \\ &= h(x) + \frac{1}{2\eta} \|x - x_t + \eta \nabla f(x_t)\|_2^2 + \text{const.}\end{aligned}$$

Define the proximal operator: $\mathbb{R}^d \rightarrow \mathbb{R}^d$.

$$\text{prox}_h(v) = \underset{x \in \mathbb{R}^d}{\operatorname{argmin}} \left\{ h(x) + \frac{1}{2} \|x - v\|_2^2 \right\}$$

从优化问题的角度看，prox算子是对函数 h 与一个二次惩罚项的联合最小化问题的解。当 η 固定时，这一问题是强凸的，且在一般的适度条件下有唯一解。

几何直观上： $\text{prox}_{\eta h}(y)$ 是在点 y 附近寻找一个点 x ，该点在尽量靠近 y 的同时，使 $h(x)$ 的值更小，从而达到折中。这正对应了不可微项的“修正”作用。

如果 $h(x)$ 是如 $\lambda \|x\|_1$ 这样的 L1 正则项，则 $\text{prox}_{\eta h}$ 的闭式解为软阈值算子，从而产生稀疏解。这也是 Proximal Gradient 被广泛用于 Lasso 等稀疏优化问题的原因。

Notice that

$$\ell(x) \leq \frac{1}{2} \|x - (x_t - \eta \nabla f(x_t))\|_2^2 + \eta h(x) + \text{const}$$

The minimizer of $\ell(x)$ is given by

$$\text{prox}_{\eta h}(x_t - \eta \cdot \nabla f(x_t))$$

Proximal gradient method: properties

Proximal gradient algorithm: for $t = 1, 2, \dots$

$$\mathbf{x}_{t+1} = \text{prox}_{\eta h}(\mathbf{x}_t - \eta \nabla f(\mathbf{x}_t)) \quad (\text{凸且L光滑})$$

- fast convergence when f is convex and L -smooth: take $\eta = 1/L$,

$$F(\mathbf{x}_t) - F^* \leq \frac{L}{2t} \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2 \quad \text{收敛速度 } O(\frac{1}{t})$$

- exponential convergence when f is μ -strongly convex $(\mu\text{-强凸指收敛})$

$$\|\mathbf{x}_t - \mathbf{x}^*\|_2^2 \leq (1 - \mu/L)^t \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2$$

- when $h(\mathbf{x}) = 0$ when $\mathbf{x} \in \mathcal{A}$ and $h(\mathbf{x}) = \infty$ otherwise, this gives the projected gradient descent for $\min_{\mathbf{x} \in \mathcal{A}} f(\mathbf{x})$:

投影梯度下降

$$\mathbf{x}_{t+1} = \mathcal{P}_{\mathcal{A}}(\mathbf{x}_t - \eta \nabla f(\mathbf{x}_t))$$

$$h(\mathbf{x}) = \begin{cases} 0 & \mathbf{x} \in \mathcal{A} \\ \infty & \mathbf{x} \notin \mathcal{A} \end{cases}$$

- Recommended reading material: Lecture 5 of the course Large-Scale Optimization for Data Science

投影回可行域 \mathcal{A}

Application to LASSO

- LASSO:

$$f(\beta) = \|\mathbf{Y} - \mathbf{X}\beta\|_2^2 \quad \text{and} \quad h(\beta) = \lambda\|\beta\|_1$$

- The proximal operator admits closed-form expression 软阈值

$$\text{prox}_h(\mathbf{v}) = \arg \min_{\beta \in \mathbb{R}^d} \left\{ \frac{1}{2} \|\beta - \mathbf{v}\|_2^2 + \lambda \|\beta\|_1 \right\} = \text{shrink}_\lambda(\mathbf{v}) \quad \text{封闭形式解}$$

where $\text{shrink}_\lambda(\cdot)$ applies entrywise shrinkage to \mathbf{v} towards zero:

$$[\text{shrink}_\lambda(\mathbf{v})]_j = \begin{cases} v_j - \lambda, & \text{if } v_j \geq \lambda, \\ v_j + \lambda, & \text{if } v_j \leq -\lambda, \\ 0, & \text{otherwise.} \end{cases}$$

何谓：当 v_j 不足，
该坐标置零 \Rightarrow 稀疏
 v_j 足够大， $v_j \pm \lambda$

- Proximal gradient algorithm for LASSO:

$$\beta_{t+1} = \text{shrink}_{\eta\lambda}(\beta_t - 2\eta \mathbf{X}^\top \mathbf{X} \beta_t + 2\eta \mathbf{X}^\top \mathbf{Y})$$

$$\beta_t - \eta \nabla f(\beta_t) = \beta_t - \eta(2\mathbf{X}^\top \mathbf{X} \beta_t - 2\mathbf{X}^\top \mathbf{Y}) = \boxed{\beta_t - 2\eta \mathbf{X}^\top \mathbf{X} \beta_t + 2\eta \mathbf{X}^\top \mathbf{Y}}$$

Solve Lasso: $\text{prox}_\lambda(v) = \arg \min_{\beta} \left\{ \frac{1}{2} \|\beta - v\|_2^2 + \lambda \|\beta\|_1 \right\}$

Lasso 问题可写为:

$$\min_{\beta \in \mathbb{R}^d} f(\beta) + h(\beta)$$

where $f(\beta) = \|Y - X\beta\|_2^2 \triangleq \text{RSS}$

$$h(\beta) = \lambda \|\beta\|_1 = \lambda \sum_{j=1}^d |\beta_j|, \quad L_1 - N$$

问题分解为一堆优化问题:

$$\ell(\beta) = \sum_{j=1}^d \underbrace{\frac{1}{2} (\beta_j - v_j)^2 + \lambda |\beta_j|}_{\ell_j(\beta_j)} = \sum_{j=1}^d \ell_j(\beta_j)$$

To minimize $\ell_j(\beta_j) = \frac{1}{2} (\beta_j - v_j)^2 + \lambda |\beta_j|$

$$= \begin{cases} \frac{1}{2} (\beta_j - v_j)^2 + \lambda \beta_j, & \text{if } \beta_j \geq 0 \\ \frac{1}{2} (\beta_j - v_j)^2 - \lambda \beta_j, & \text{if } \beta_j \leq 0 \end{cases}$$

$$= \begin{cases} \frac{1}{2} \beta_j^2 + (\lambda - v_j) \beta_j + \frac{1}{2} v_j^2, & \text{if } \beta_j \geq 0 \\ \frac{1}{2} \beta_j^2 - (\lambda + v_j) \beta_j + \frac{1}{2} v_j^2, & \text{if } \beta_j \leq 0 \end{cases}$$

Within $\beta_j \geq 0$, $\ell_j(\beta_j)$ is minimized at $\begin{cases} \lambda - v_j, & \text{if } \lambda - v_j \geq 0 \\ 0, & \text{if } \lambda - v_j \leq 0 \end{cases}$

Within $\beta_j \leq 0$, $\ell_j(\beta_j)$ is minimized at $\begin{cases} 0, & \text{if } \lambda + v_j \geq 0 \\ \lambda + v_j, & \text{if } \lambda + v_j \leq 0 \end{cases}$

$\hat{\beta}_j = \arg \min_{\beta} \ell(\beta_j)$

$\Rightarrow \hat{\beta}_j = \begin{cases} \lambda - v_j, & v_j \geq \lambda \\ 0, & -\lambda \leq v_j \leq \lambda \\ \lambda + v_j, & v_j \leq -\lambda \end{cases}$

$= S_v(v_j) = \text{sign}(v_j) \max(|v_j| - \lambda, 0)$

对非光滑凸函数, 定义 次梯度 (Subgradient)

for convex function f : define g is a subgradient at point x , if for $\forall y$.

$$f(y) \geq f(x) + g^T(y - x)$$

Example: $f(x) = |x|$. $\partial f(0) = [-1, 1]$

check: 0 is minimizer of $f \Leftrightarrow 0 \in \partial f(x)$

Pf: " \Rightarrow " if x is minimizer, $\forall y \in \mathbb{R}$

$$f(y) \geq f(x) = f(x) + \langle 0, y - x \rangle \Rightarrow 0 \in \partial f(x)$$

" \Leftarrow " if $0 \in \partial f(x)$, $\forall y \in \mathbb{R}$, $f(y) \geq f(x) + \langle 0, y - x \rangle = f(x)$ \square

If $\ell(\beta) = \sum_{j=1}^d \frac{1}{2} \beta_j^2 + \lambda \beta_j$

$$= \sum_{j=1}^d \frac{1}{2} \beta_j^2 + \lambda \beta_j = \sum_{j=1}^d \frac{1}{2} \beta_j^2 - \left(\lambda - \frac{v_j}{2}\right) \beta_j + \frac{1}{2} v_j^2$$

$$= \sum_{j=1}^d \frac{1}{2} \beta_j^2 - \left(\lambda + \frac{v_j}{2}\right) \beta_j + \frac{1}{2} v_j^2$$

Statistical properties of LASSO

Setup

LASSO:

$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^d} \left\{ \frac{1}{2} \|\mathbf{Y} - \mathbf{X}\beta\|_2^2 + \lambda \|\beta\|_1 \right\},$$

- Independent, sub-Gaussian noise $\|\varepsilon_i\|_{\psi_2} \leq \sigma$ 独立次高斯噪声
- Sparsity: $n \gg s \log d$ 稀疏性, 记 $S = \{j : \beta_j^* \neq 0\}$, $|S| = s \ll d$

- Theory-informed tuning parameter selection:

为了良好的特征选择能力和估计精度

$$\lambda \asymp \sigma \sqrt{n \log d}$$

真实参数支撑集

- Question:

一致性

- Does LASSO recover the support of β^* ? model selection consistency
- Does LASSO provide reliable estimate for β^* ?

$$\sqrt{\frac{s \log d}{n}} \text{ 量级}$$

Optimality condition

The optimality condition for unconstrained convex optimization

$$\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x})$$

- if f is smooth: $\nabla f(\hat{\mathbf{x}}) = \mathbf{0}$
- in general (when f might not be smooth): $\mathbf{0} \in \partial f(\hat{\mathbf{x}})$

Here $\boxed{\partial f(\mathbf{x}) \subseteq \mathbb{R}^d}$ is the **subgradient** of the convex function f at \mathbf{x} :

$$\mathbf{g} \in \partial f(\mathbf{x}) \iff f(\mathbf{y}) \geq f(\mathbf{x}) + \mathbf{g}^\top (\mathbf{y} - \mathbf{x}) \quad \text{for all } \mathbf{y} \in \mathbb{R}^d$$



Check (in homework):

- if f is smooth at \mathbf{x} : $\partial f(\mathbf{x}) = \{\nabla f(\mathbf{x})\}$
- the optimality condition for LASSO is: for each $1 \leq j \leq d$

$$[\mathbf{X}^\top (\mathbf{Y} - \mathbf{X}^\top \hat{\boldsymbol{\beta}})]_j \quad \begin{cases} = \lambda \cdot \text{sign}(\hat{\beta}_j) & \text{if } \hat{\beta}_j \neq 0 \\ \in [-\lambda, \lambda] & \text{if } \hat{\beta}_j = 0 \end{cases}$$