

Sparse Solutions to Least-Squares Problems Using the LASSO

Objectives

- motivate search for sparse solutions
- introduce ℓ_1 -norm regularization (LASSO)
- overview attributes of ℓ_1 -regularization

Sparse classifiers/models give insight 2

$(\underline{x}_i, d_i), i=1, \dots, N$

→ 框疏模型就是模型中的绝大部分参数都是0 (意味着只有少部分特征是重要的)

features, labels

$$\underline{x}_i^\top \underline{w} \approx d_i$$

$$A \underline{w} = \begin{bmatrix} \underline{a}_1 & \underline{a}_2 & \cdots & \underline{a}_m \end{bmatrix}_{N \times m} \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_m \end{bmatrix} = \sum_{i=1}^m w_i \underline{a}_i$$

\underline{a}_l : lth feature component

对所有特征的加权平均

Suppose $w_e \approx 0 \Rightarrow \underline{a}_e$ is unimportant

If a small number of w_i are nonzero, only those few features matter! \underline{w} is sparse

$$\|\underline{w}\|_0 = \sum_{i=1}^m \mathbb{1}_{\{w_i \neq 0\}}$$

非零 w 的个数 (number of nonzero elements)

$\|\cdot\|_0$ "norm"

$$\|\underline{a} \underline{w}\|_0 \neq \|\underline{a}\| \|\underline{w}\|_0$$

其实不满足
数学上对
范数的定义

$$\text{Consider } \min_{\underline{w}} \|\underline{w}\|_0 \text{ s.t. } \|A \underline{w} - \underline{d}\|_2 \leq \varepsilon$$

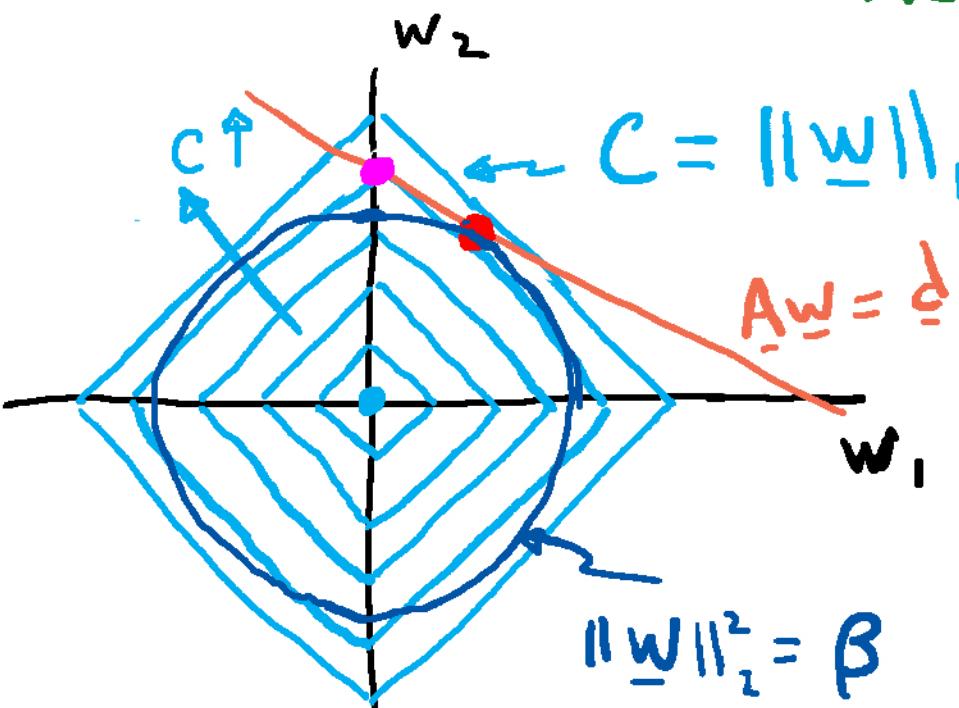
The more sparse, the better

non convex-intractable

Convex relaxation gives tractable problem

$$\min_{\underline{w}} \|\underline{w}\|_1 \text{ s.t. } \|\underline{A}\underline{w} - \underline{d}\|_2^2 \leq \Sigma \quad \text{LASSO: Least Absolute Selection + Shrinkage Operator}$$

This solution is tractable



$$\min \|\underline{w}\|_2^2 \text{ s.t. } \underline{A}\underline{w} = \underline{d}$$

circular $\|\underline{w}\|_2^2 \Rightarrow$ non sparse solutions

\downarrow w_1 情况下 $Aw=d$ 的直线和红点交于 (w_1, w_2)
(w_1 和 w_2 都不 = 0)

$$C = \|\underline{w}\|_1 = \sum_{i=1}^m |w_i| : |w_1| + |w_2| = C$$

1st quad $w_1 + w_2 = C$

$\min \|\underline{w}\|_1 \text{ s.t. } \underline{A}\underline{w} = \underline{d}$
"corners" on $\|\underline{w}\|_1 \Rightarrow$ sparse solns

$\begin{matrix} w_1 \\ \parallel \\ (0, w_2) \end{matrix}$

LASSO is a regularized least-squares problem 4

$\min_{\underline{w}} \|\underline{w}\|_1$, s.t. $\|\underline{A}\underline{w} - \underline{d}\|_2^2 \leq \varepsilon$ is equivalent to

$\min_{\underline{w}} \|\underline{A}\underline{w} - \underline{d}\|_2^2 + \lambda \|\underline{w}\|_1$, for some λ, ε

Note: $\min_{\underline{w}} \|\underline{w}\|_1 + \frac{1}{\lambda} \|\underline{A}\underline{w} - \underline{d}\|_2^2$ If we factor out the λ ,
this is the same with the above line

LASSO

$$\underline{w}_L = \arg \min_{\underline{w}} \|\underline{A}\underline{w} - \underline{d}\|_2^2 + \lambda \|\underline{w}\|_1$$

Sparse $\underline{w}_{L_{\text{lasso}}}$

can have small model error

$$\underline{w}_{\text{opt}} - \underline{w}_L$$

iterative solution

Ridge Regression

$$\underline{w}_R = \arg \min_{\underline{w}} \|\underline{A}\underline{w} - \underline{d}\|_2^2 + \lambda \|\underline{w}\|_2^2$$

non sparse $\underline{w}_{R_{\text{ridge}}}$

great prediction error

$$\|\underline{A}\underline{w}_{\text{opt}} - \underline{A}\underline{w}_R\|_2^2$$
_(optimal) _(ridge regression solution)

can solve in closed form

LASSO may be used for model/feature selection 5

$$\underline{w}_L = \arg \min_{\underline{w}} \left\| \underline{A} \underline{w} - \underline{d} \right\|_2^2 + \lambda \|\underline{w}\|_1$$

Model Error
Sparsity penalty

$$S_L = \{i : [\underline{w}_L]_i \neq 0\}$$

selected features

只使用不重要的特征
的权重为0，我们选择出那些权重不为0的特征

$$\underline{A} \underline{w}_L = \sum_{i=1}^m \underline{a}_i [\underline{w}_L]_i = \sum_{i \in S_L} \underline{a}_i [\underline{w}_L]_i$$

Debiasing

$$\underline{A}_L = \{\underline{a}_i : i \in S_L\}$$

$$\hat{\underline{w}}_L = \arg \min_{\underline{w}} \left\| \underline{A}_L \underline{w} - \underline{d} \right\|_2^2 = (\underline{A}_L^\top \underline{A}_L)^{-1} \underline{A}_L^\top \underline{d}$$

avoids shrinkage due to $\|\underline{w}\|_1$

- ① 只用 LASSO 选 feature
- ② 模型全训 → 组成一个更小的 A 做标准 LS

Copyright 2019
Barry Van Veen