

Other regularizations, sparsity and the Lasso

MATH-412 - Statistical Machine Learning

A diverse set of regularization approaches

Regularizers are not necessarily quadratic and not necessarily norms :

$$\min_{f \in S} \widehat{\mathcal{R}}(f) + \lambda \Omega(f) \quad \text{with e.g.} \quad \Omega(f) = \int (f''(\mathbf{x}))^2 d\mathbf{x} \quad \text{or} \quad \int |f'(\mathbf{x})| d\mathbf{x}$$

It is possible to couple the regularization of different tasks :

$$\min_{f_1, f_2, \dots, f_K \in S} \sum_k \widehat{\mathcal{R}}_{(k)}(f_k) + \lambda \sum_k \|f_k - \bar{f}\|^2 + \mu \|\bar{f}\|^2$$

Even when the predictor has a linear parameterisation, there are various options

$$\min_{\mathbf{w} \in \mathbb{R}^p} \widehat{\mathcal{R}}(\mathbf{w}) + \lambda \Omega(\mathbf{w}) \quad \text{with} \quad \Omega(\mathbf{w}) = \dots$$

$$\|\mathbf{w}\|_q \quad , \quad \|\mathbf{w}\|_1 \quad , \quad \|\mathbf{w}\|_1 + \eta \|\mathbf{w}\|_2^2 \quad , \quad \sum_{j=1}^{p-1} (w_{j+1} - w_j)^2 \quad , \quad \sum_{j=1}^{p-1} |w_{j+1} - w_j| \quad , \quad \text{etc.}$$

Best subset selection

The number of effectively used variables is often denoted

$$\|\mathbf{w}\|_0 := \#\{j \mid w_j \neq 0\} = \sum_{j=1}^p 1_{\{w_j \neq 0\}}$$

Best subset selection formulation

$$\min_{\mathbf{w} \in \mathbb{R}^p} \hat{\mathcal{R}}(\mathbf{w}) + \lambda \|\mathbf{w}\|_0$$

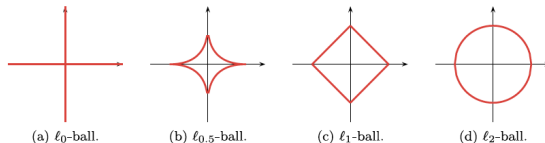
- Compromise between fitting and $\#$ of variables in the model
- The problem is NP-hard to solve in general
- Can be solved by exhaustive search amongst 2^p models for p small

Lasso (Least Absolute Shrinkage and Selection Operator)

$$\min_{\mathbf{w} \in \mathbb{R}^p} \hat{\mathcal{R}}(\mathbf{w}) + \lambda \|\mathbf{w}\|_1$$

- No closed form solution
- Convex but non-differentiable optimization problem
- Can nonetheless be solved by efficient algorithms

The general approach extends to $q < 1$ with quasi-norms but then the problem **is not convex anymore**.



Lasso regression : Constrained vs regularized problem

$$\min_{\mathbf{w} \in \mathbb{R}^p} \frac{1}{2n} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_1$$

vs

$$\min_{\mathbf{w} \in \mathbb{R}^p} \frac{1}{2n} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 \quad \text{s.t.} \quad \|\mathbf{w}\|_1 \leq C$$

$$\min_{\mathbf{w} \in \mathbb{R}^p} f(\mathbf{w}) + \lambda g(\mathbf{w})$$

vs

$$\min_{\mathbf{w} \in \mathbb{R}^p} f(\mathbf{w}) \quad \text{s.t.} \quad g(\mathbf{w}) \leq C$$

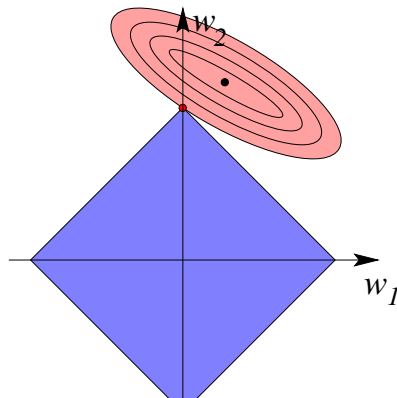
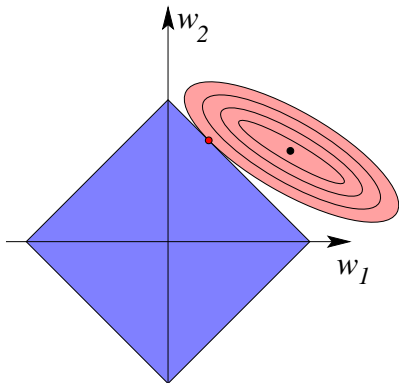
Proposition

If f and g are convex, then for any value of λ there is a value of C such that both problems have the same solution and vice versa.

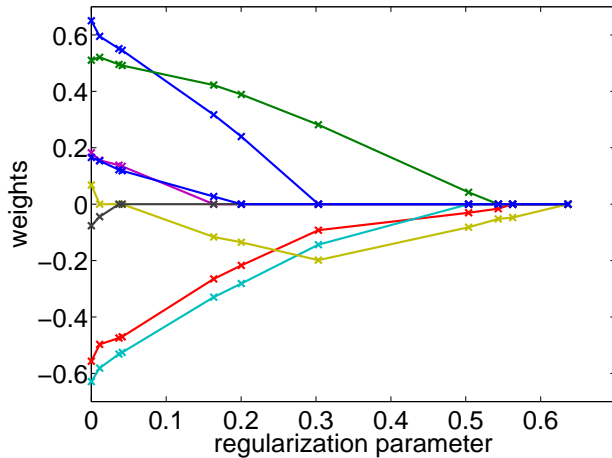
Geometric intuition for the Lasso

Consider the constrained problem

$$\min_{\mathbf{w} \in \mathbb{R}^p} \frac{1}{2n} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 \quad \text{s.t.} \quad \|\mathbf{w}\|_1 \leq C$$



Lasso regression has piecewise linear paths



Lasso regression with orthogonal design

- Assume $\frac{1}{n}\mathbf{X}^\top\mathbf{X} = \mathbf{I}$
- Then solving the Lasso is equivalent to solving

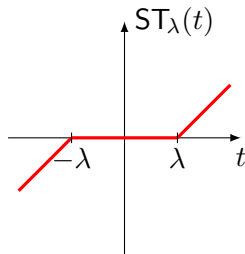
$$\min_w \frac{1}{2}\|\hat{\mathbf{c}} - \mathbf{w}\|_2^2 + \lambda\|\mathbf{w}\|_1 \quad \text{for} \quad \hat{\mathbf{c}} = \frac{1}{n}\mathbf{X}^\top\mathbf{y}$$

- $\hat{c}_j = \frac{1}{n}\mathbf{y}^\top\mathbf{x}^{(j)}$
- Equivalent to solve $\forall j$

$$\begin{aligned}\hat{w}_j &= \arg \min_{v \in \mathbb{R}} \frac{1}{2}v^2 - v\hat{c}_j + \lambda|v| \\ &= \text{ST}_\lambda(\hat{c}_j)\end{aligned}$$

with the *soft-thresholding* operator :

$$\text{ST}_\lambda(t) := (|t| - \lambda)_+ \text{ sign}(t)$$



Best subset selection with orthogonal design

- Assume $\frac{1}{n}\mathbf{X}^\top\mathbf{X} = \mathbf{I}$
- Then solving the Lasso is equivalent to solving

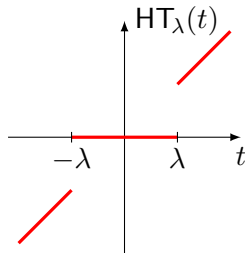
$$\min_w \frac{1}{2}\|\hat{\mathbf{c}} - \mathbf{w}\|_2^2 + \lambda\|\mathbf{w}\|_0 \quad \text{for} \quad \hat{\mathbf{c}} = \frac{1}{n}\mathbf{X}^\top\mathbf{y}$$

- $\hat{c}_j = \frac{1}{n}\mathbf{y}^\top\mathbf{x}^{(j)}$
- Equivalent to solve $\forall j$

$$\begin{aligned}\hat{w}_j &= \arg\min_{v \in \mathbb{R}} \frac{1}{2}v^2 - v\hat{c}_j + \lambda 1_{\{v \neq 0\}} \\ &= \text{HT}_\lambda(\hat{c}_j)\end{aligned}$$

with the *hard-thresholding* operator :

$$\text{HT}_\lambda(t) := t 1_{\{|t| > \lambda\}}$$



Tackling the ℓ_0 constrained problem for p large...

$$\min_{\mathbf{w} \in \mathbb{R}^p} \frac{1}{2n} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 \quad \text{s.t.} \quad \|\mathbf{w}\|_0 \leq k$$

- The problem is NP-hard : what if p is large ?

Greedy methods

Principle : \mathbf{w} is estimated by increasing the support greedily. At each iteration

- 1 **Selection step** : A new coordinate is included in the support of \mathbf{w}
- 2 **Fitting step** : The new coefficient and possibly old ones are re-optimized

Forward selection (regression)

Initialization :

- $\hat{S} = \emptyset$ (estimate of support)

Repeat :

① **Selection Step :**

- $j \leftarrow \arg \min_{j'} \min_{\mathbf{w}_{\hat{S} \cup \{j'\}}} \|\mathbf{y} - \mathbf{X}_{\hat{S} \cup \{j'\}} \mathbf{w}_{\hat{S} \cup \{j'\}}\|_2^2, \quad \hat{S} \leftarrow \hat{S} \cup \{j\}$

② **Fitting Step :**

- $\hat{\mathbf{w}}_{\hat{S}} \leftarrow \arg \min_{\mathbf{w}_{\hat{S}}} \|\mathbf{y} - \mathbf{X}_{\hat{S}} \mathbf{w}_{\hat{S}}\|_2^2$

Backward selection :

- Symmetric by removing variables one by one
- Not recommended if the number of variables is large, because starting from an overfitted situation.

Orthogonal Matching Pursuit (regression)

Initialization :

- $\hat{S} = \emptyset$ (estimate of support)
- $\mathbf{r} \leftarrow \mathbf{y}$ (residuals)

Repeat :

① Selection Step :

- $j \leftarrow \arg \max_{j'} |\langle \mathbf{x}^{(j')}, \mathbf{r} \rangle|, \quad \hat{S} \leftarrow \hat{S} \cup \{j\}$

② Fitting Step :

- $\hat{\mathbf{w}}_{\hat{S}} \leftarrow \arg \min_{\mathbf{w}_{\hat{S}}} \|\mathbf{y} - \mathbf{X}_{\hat{S}} \mathbf{w}_{\hat{S}}\|_2^2 \quad \mathbf{r} \leftarrow \mathbf{y} - \mathbf{X}_{\hat{S}} \hat{\mathbf{w}}_{\hat{S}}$

Comparing Lasso and other strategies for linear regression

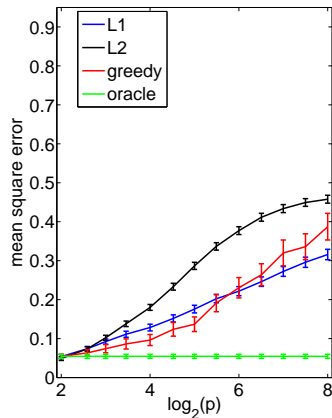
Comparing :

$$\begin{aligned}\text{Ridge regression : } & \min_{\mathbf{w} \in \mathbb{R}^p} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \frac{\lambda}{2} \|\mathbf{w}\|_2^2 \\ \text{Lasso : } & \min_{\mathbf{w} \in \mathbb{R}^p} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_1 \\ \text{OMP/FS : } & \min_{\mathbf{w} \in \mathbb{R}^p} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_0\end{aligned}$$

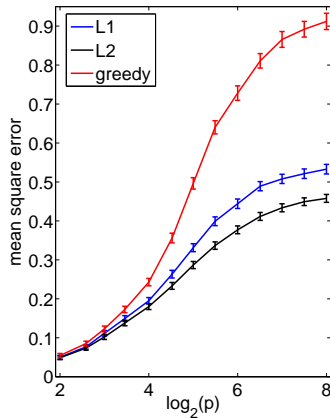
- Each method builds a path of solutions from 0 to ordinary least-squares solution
- Regularization parameters selected on the test set

Simulation results

- \mathbf{X} = i.i.d. Gaussian design, $n = 64$, $p \in [2, 256]$,
- $\mathbf{y} = \mathbf{X}\mathbf{w}^* + \boldsymbol{\varepsilon}$, $\|\mathbf{w}^*\|_0 = 4$, $w_j^* \in \{-1, 0, 1\}$, $\sigma^2 = 1$.



Sparse



Rotated (non sparse)

Note ℓ_1 stability
to non-sparsity