

The grouping effect

- We consider the generic penalization method

$$\hat{\beta} = \arg \min_{\beta} |\mathbf{y} - \mathbf{X}\beta|^2 + \lambda J(\beta) \quad (5)$$

where $J(\cdot)$ is positive valued for $\beta \neq 0$

- **Lemma1** Assume that $\mathbf{x}_i = \mathbf{x}_j, i, j \in \{1, \dots, p\}$.
 - (a) If $J(\cdot)$ is strictly convex, then $\hat{\beta}_i = \hat{\beta}_j, \forall \lambda > 0$.
 - (b) If $J(\beta) = |\beta|_1$, then $\hat{\beta}_i \hat{\beta}_j \geq 0$ and $\hat{\beta}^*$ is another minimizer of equation (5), where

$$\hat{\beta}_k^* = \begin{cases} \hat{\beta}_k & \text{if } k \neq i \text{ and } k \neq j, \\ (\hat{\beta}_i + \hat{\beta}_j) \cdot s & \text{if } k = i, \\ (\hat{\beta}_i + \hat{\beta}_j) \cdot (1 - s) & \text{if } k = j. \end{cases}$$

for any $s \in [0, 1]$

Proof of the lemma1

► Part(a)

Fix $\lambda > 0$. If $\hat{\beta}_i \neq \hat{\beta}_j$, let us consider $\hat{\beta}^*$ as follows:

$$\hat{\beta}_k^* = \begin{cases} \hat{\beta}_k & \text{if } k \neq i \text{ and } k \neq j, \\ \frac{1}{2}(\hat{\beta}_i + \hat{\beta}_j) & \text{if } k = i \text{ and } k = j \end{cases}$$

Because $\mathbf{x}_i = \mathbf{x}_j$, it is obvious that $\mathbf{X}\hat{\beta}^* = \mathbf{X}\hat{\beta}$; thus $|\mathbf{y} - \mathbf{X}\hat{\beta}^*|^2 = |\mathbf{y} - \mathbf{X}\hat{\beta}|^2$. However, $J(\cdot)$ is strictly convex, so we have $J(\hat{\beta}^*) < J(\hat{\beta})$. Therefore, $\hat{\beta}$ cannot be the minimizer of equation (5), which is a contradiction. So we must have $\hat{\beta}_i = \hat{\beta}_j$.

Lemma 1 (Part b)

Statement: Assume $\mathbf{x}_i = \mathbf{x}_j$, $i, j \in \{1, \dots, p\}$. If $J(\boldsymbol{\beta}) = \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 + \lambda\|\boldsymbol{\beta}\|_1$, then:

- ▶ $\hat{\beta}_i \cdot \hat{\beta}_j > 0$,
- ▶ $\hat{\boldsymbol{\beta}}^*$ is another minimizer of $J(\boldsymbol{\beta})$, where:

$$\hat{\beta}_k^* = \begin{cases} \hat{\beta}_k, & \text{if } k \neq i \text{ and } k \neq j, \\ (\hat{\beta}_i + \hat{\beta}_j) \cdot s, & \text{if } k = i, \\ (\hat{\beta}_i + \hat{\beta}_j) \cdot (1 - s), & \text{if } k = j, \end{cases}$$

for any $s \in [0, 1]$.

Proof Setup (Part b of Lemma 1)

- ▶ Assume $\hat{\beta}_i \cdot \hat{\beta}_j < 0$, i.e., $\hat{\beta}_i$ and $\hat{\beta}_j$ have opposite signs.
- ▶ Define $\hat{\beta}^*$ as:

$$\hat{\beta}_k^* = \begin{cases} \hat{\beta}_k, & \text{if } k \neq i \text{ and } k \neq j, \\ (\hat{\beta}_i + \hat{\beta}_j) \cdot s, & \text{if } k = i, \\ (\hat{\beta}_i + \hat{\beta}_j) \cdot (1 - s), & \text{if } k = j. \end{cases}$$

- ▶ We aim to compare the objective values $J(\hat{\beta})$ and $J(\hat{\beta}^*)$.

Objective Function Comparison

1. Residual Term:

- ▶ Since $\mathbf{x}_i = \mathbf{x}_j$, the contribution of $\hat{\beta}_i$ and $\hat{\beta}_j$ to the residual term is unchanged:

$$\|\mathbf{y} - \mathbf{X}\hat{\beta}^*\|^2 = \|\mathbf{y} - \mathbf{X}\hat{\beta}\|^2.$$

2. L_1 -Norm Penalty Term:

- ▶ The L_1 -norm for $\hat{\beta}^*$ is:

$$\|\hat{\beta}^*\|_1 = \sum_{k \neq i,j} |\hat{\beta}_k| + |\hat{\beta}_i + \hat{\beta}_j|.$$

- ▶ Since $\hat{\beta}_i \cdot \hat{\beta}_j < 0$, we have:

$$|\hat{\beta}_i + \hat{\beta}_j| < |\hat{\beta}_i| + |\hat{\beta}_j|.$$

- ▶ Hence:

$$\|\hat{\beta}^*\|_1 < \|\hat{\beta}\|_1.$$

Conclusion of Proof (Part b)

- ▶ The total objective function for $\hat{\beta}^*$ is:

$$J(\hat{\beta}^*) = \|\mathbf{y} - \mathbf{X}\hat{\beta}^*\|^2 + \lambda \|\hat{\beta}^*\|_1.$$

- ▶ Since:

$$\|\mathbf{y} - \mathbf{X}\hat{\beta}^*\|^2 = \|\mathbf{y} - \mathbf{X}\hat{\beta}\|^2,$$

and:

$$\|\hat{\beta}^*\|_1 < \|\hat{\beta}\|_1,$$

it follows that:

$$J(\hat{\beta}^*) < J(\hat{\beta}).$$

- ▶ This contradicts the assumption that $\hat{\beta}$ minimizes $J(\beta)$.
- ▶ Therefore, $\hat{\beta}_i \cdot \hat{\beta}_j \geq 0$. Moreover, if $\hat{\beta}_i \cdot \hat{\beta}_j > 0$, the modified solution $\hat{\beta}^*$ is also a minimizer.

Conclusion of Lemma 1

Lemma 1 Summary:

- ▶ Assume $\mathbf{x}_i = \mathbf{x}_j$, where $i, j \in \{1, \dots, p\}$.
- ▶ The conclusions are:

1. If $J(\cdot)$ is strictly convex, then:

$$\hat{\beta}_i = \hat{\beta}_j, \quad \forall \lambda > 0.$$

2. If $J(\beta) = \|\beta\|_1$, then:

- ▶ $\hat{\beta}_i \cdot \hat{\beta}_j > 0$, meaning $\hat{\beta}_i$ and $\hat{\beta}_j$ have the same sign.
- ▶ Any convex combination of $\hat{\beta}_i$ and $\hat{\beta}_j$, defined as:

$$\hat{\beta}_k^* = \begin{cases} \hat{\beta}_k, & \text{if } k \neq i \text{ and } k \neq j, \\ (\hat{\beta}_i + \hat{\beta}_j) \cdot s, & \text{if } k = i, \\ (\hat{\beta}_i + \hat{\beta}_j) \cdot (1 - s), & \text{if } k = j, \end{cases}$$

for $s \in [0, 1]$, is also a minimizer of $J(\beta)$.

Theorem

- ▶ Lemma shows a clear distinction between strictly convex penalty functions and the lasso penalty. Strict convexity guarantees the grouping effect in the extreme situation with identical predictors. In contrast the lasso does not even have a unique solution.
- ▶ *Theorem 1* Given data (\mathbf{y}, \mathbf{X}) and parameters (λ_1, λ_2) , the response \mathbf{y} is centered and the predictors \mathbf{X} are standardized. Let $\hat{\beta}(\lambda_1, \lambda_2)$ be the naive elastic net estimate. Suppose that $\hat{\beta}_i(\lambda_1, \lambda_2) \hat{\beta}_j(\lambda_1, \lambda_2) > 0$ define

$$D_{\lambda_1, \lambda_2}(i, j) = \frac{1}{|\mathbf{y}|_1} |\hat{\beta}_i(\lambda_1, \lambda_2) - \hat{\beta}_j(\lambda_1, \lambda_2)|$$

then

$$D_{\lambda_1, \lambda_2}(i, j) \leq \frac{1}{\lambda_2} \sqrt{\{2(1 - \rho)\}}$$

, where $\rho = \mathbf{x}_i^T \mathbf{x}_j$, the sample correlation.

Proof of the Theorem1

If $\hat{\beta}_i(\lambda_1, \lambda_2)\hat{\beta}_j(\lambda_1, \lambda_2) > 0$, then both $\hat{\beta}_i(\lambda_1, \lambda_2)$ and $\hat{\beta}_j(\lambda_1, \lambda_2)$ are non-zero, and we have $\text{sign}\{\hat{\beta}_i(\lambda_1, \lambda_2)\} = \text{sign}\{\hat{\beta}_j(\lambda_1, \lambda_2)\}$. Because of equation (3), $\hat{\beta}(\lambda_1, \lambda_2)$ satisfies

$$\left. \frac{\partial L(\lambda_1, \lambda_2, \beta)}{\partial \beta_k} \right|_{\beta=\hat{\beta}(\lambda_1, \lambda_2)} = 0 \quad \text{if } \hat{\beta}_k(\lambda_1, \lambda_2) \neq 0.$$

Hence we have

$$-2\mathbf{x}_i^\top \{\mathbf{y} - \mathbf{X}\hat{\beta}(\lambda_1, \lambda_2)\} + \lambda_1 \text{sgn}\{\hat{\beta}_i(\lambda_1, \lambda_2)\} + 2\lambda_2 \hat{\beta}_i(\lambda_1, \lambda_2) = 0, \quad (6)$$

$$-2\mathbf{x}_j^\top \{\mathbf{y} - \mathbf{X}\hat{\beta}(\lambda_1, \lambda_2)\} + \lambda_1 \text{sgn}\{\hat{\beta}_j(\lambda_1, \lambda_2)\} + 2\lambda_2 \hat{\beta}_j(\lambda_1, \lambda_2) = 0. \quad (7)$$

Subtracting equation (6) from equation (7) gives

$$(\mathbf{x}_j^\top - \mathbf{x}_i^\top)\{\mathbf{y} - \mathbf{X}\hat{\beta}(\lambda_1, \lambda_2)\} + \lambda_2\{\hat{\beta}_i(\lambda_1, \lambda_2) - \hat{\beta}_j(\lambda_1, \lambda_2)\} = 0,$$

which is equivalent to

$$\hat{\beta}_i(\lambda_1, \lambda_2) - \hat{\beta}_j(\lambda_1, \lambda_2) = \frac{1}{\lambda_2}(\mathbf{x}_i^\top - \mathbf{x}_j^\top)\hat{\mathbf{r}}(\lambda_1, \lambda_2), \quad (8)$$

where $\hat{\mathbf{r}}(\lambda_1, \lambda_2) = \mathbf{y} - \mathbf{X}\hat{\beta}(\lambda_1, \lambda_2)$ is the residual vector,

Proof of the Theorem1

Since \mathbf{X} is standardized,

$$|\mathbf{x}_i - \mathbf{x}_j|^2 = 2(1 - \rho), \quad \rho = \mathbf{x}_i^\top \mathbf{x}_j.$$

By equation (3), we must have

$$L\{\lambda_1, \lambda_2, \hat{\beta}(\lambda_1, \lambda_2)\} \leq L(\lambda_1, \lambda_2, \beta = 0),$$

i.e.,

$$|\hat{r}(\lambda_1, \lambda_2)|^2 + \lambda_2 |\hat{\beta}(\lambda_1, \lambda_2)|^2 + \lambda_1 |\hat{\beta}(\lambda_1, \lambda_2)|_1 \leq |\mathbf{y}|^2.$$

So $|\hat{r}(\lambda_1, \lambda_2)| \leq |\mathbf{y}|$. Then equation (8) implies that

$$D_{\lambda_1, \lambda_2}(i, j) \leq \frac{1}{\lambda_2} \frac{|\hat{r}(\lambda_1, \lambda_2)|}{|\mathbf{y}|} |\mathbf{x}_i - \mathbf{x}_j| \leq \frac{1}{\lambda_2} \sqrt{2(1 - \rho)}.$$

Conclusion of Theorem 1

Theorem 1 Summary:

- ▶ Given standardized predictors \mathbf{X} and centered response \mathbf{y} , let $\hat{\beta}(\lambda_1, \lambda_2)$ be the naive elastic net estimate.
- ▶ Suppose $\hat{\beta}_i(\lambda_1, \lambda_2) \cdot \hat{\beta}_j(\lambda_1, \lambda_2) > 0$, for predictors i and j .
- ▶ Define:

$$D_{\lambda_1, \lambda_2}(i, j) = \frac{1}{\|\mathbf{y}\|_1} \cdot |\hat{\beta}_i(\lambda_1, \lambda_2) - \hat{\beta}_j(\lambda_1, \lambda_2)|.$$

- ▶ Then:

$$D_{\lambda_1, \lambda_2}(i, j) \leq \frac{1}{\lambda_2} \sqrt{2(1 - \rho)},$$

where $\rho = \mathbf{x}_i^\top \mathbf{x}_j$ is the sample correlation between predictors i and j .

Conclusion:

- ▶ The difference between the coefficients of highly correlated predictors is bounded by $D_{\lambda_1, \lambda_2}(i, j)$.
- ▶ This demonstrates the **grouping effect** of the naive elastic net, ensuring similar coefficients for highly correlated predictors.

Why Lasso Does Not Exhibit the Grouping Effect

1. Independent Treatment of Predictors:

The L_1 penalty treats each predictor independently, without considering correlations between them. As a result, lasso arbitrarily selects one predictor from a group of highly correlated variables and shrinks the rest to zero.

2. Geometric Explanation:

- ▶ Lasso's constraint region is diamond-shaped, which encourages sparsity by forcing some coefficients to zero.
- ▶ For highly correlated predictors, the optimal solution often lies on a vertex of the diamond, where only one coefficient is nonzero. This breaks the grouping effect.

3. Example from Tibshirani (1996):

- ▶ Consider a simple linear model with two predictors ($p = 2$)

$$|\hat{\beta}_1 - \hat{\beta}_2| = |\cos(\theta)|$$

where, θ : Angle between \mathbf{y} and the difference vector $\mathbf{x}_1 - \mathbf{x}_2$

- ▶ As $\rho = \text{corr}(\mathbf{x}_1, \mathbf{x}_2) \rightarrow 1$, $\cos(\theta)$ does not necessarily vanish. This means the coefficients $\hat{\beta}_1$ and $\hat{\beta}_2$ may still differ significantly, even for highly correlated predictors.