

Stat 435 Intro to Statistical Machine Learning

Week 6: Additional exercises

Richard Li

May 9, 2017

How do they compare

$$\begin{aligned}\|\vec{u}\|_2^2 &= \sum \kappa_i^2 \\ &= \vec{u}^T \vec{u}\end{aligned}$$

Assuming all \mathbf{X} and \mathbf{Y} are centered,

$$\beta_1 = \operatorname{argmin} \|\mathbf{Y} - \mathbf{X}\beta\|_2^2$$

$$\beta_2 = \operatorname{argmin} \|\mathbf{Y} - \mathbf{X}\beta\|_2^2 + \lambda\|\beta\|_2^2$$

$$\beta_3 = \operatorname{argmin} \|\mathbf{Y} - a\mathbf{X}\beta\|_2^2$$

$$\beta_4 = \operatorname{argmin} \|\mathbf{Y} - a\mathbf{X}\beta\|_2^2 + \lambda\|\beta\|_2^2$$

$$\beta_5 = \operatorname{argmin} \|\mathbf{Y} - a\mathbf{X}\beta\|_2^2 + a^2\lambda\|\beta\|_2^2$$

$$\{\beta_0, \beta_6\} = \operatorname{argmin} \|(\mathbf{Y} + 1) - \beta_0 - (\mathbf{X} + 2)\beta\|_2^2$$

$$\{\beta_0, \beta_7\} = \operatorname{argmin} \|(\mathbf{Y} + 1) - \beta_0 - (\mathbf{X} + 2)\beta\|_2^2 + \lambda\|\beta\|_2^2$$

And what about the resulted RSS?

Least square: shifting by a constant

$$\beta_1 = \operatorname{argmin} \|\mathbf{Y} - \mathbf{X}\beta\|_2^2$$

$$\{\beta_0, \beta_6\} = \operatorname{argmin} \|(\mathbf{Y} + 1) - \beta_0 - (\mathbf{X} + 2)\beta\|_2^2$$

- $\beta_6 = \beta_1$
- We can calculate β_0 from β_6 , $\bar{\mathbf{X}}$, and $\bar{\mathbf{Y}}$
- RSS is also the same

Least square: changing scales

$$\beta_1 = \operatorname{argmin} \|\mathbf{Y} - \mathbf{X}\beta\|_2^2$$

$$\beta_3 = \operatorname{argmin} \|\mathbf{Y} - a\mathbf{X}\beta\|_2^2$$

- $\beta_3 = \beta_1/a$
- since $\mathbf{Y} - \mathbf{X}\beta_1 = \mathbf{Y} - a\mathbf{X}(\beta_1/a)$
- RSS is also the same

Ridge: shifting by a constant

$$\beta_2 = \operatorname{argmin} \|\mathbf{Y} - \mathbf{X}\beta\|_2^2 + \lambda\|\beta\|_2^2$$

$$\{\beta_0, \beta_7\} = \operatorname{argmin} \|(\mathbf{Y} + 1) - \beta_0 - (\mathbf{X} + 2)\beta\|_2^2 + \lambda\|\beta\|_2^2$$

- $\beta_7 = \beta_2$
- Notice β_0 is not penalized, so similar to the argument before
- We can calculate β_0 from β_7 , $\bar{\mathbf{X}}$, and $\bar{\mathbf{Y}}$
- RSS is also the same

Ridge: compare with least square

$$\beta_1 = \operatorname{argmin} \|\mathbf{Y} - \mathbf{X}\beta\|_2^2$$

$$\beta_2 = \operatorname{argmin} \|\mathbf{Y} - \mathbf{X}\beta\|_2^2 + \lambda \|\beta\|_2^2$$

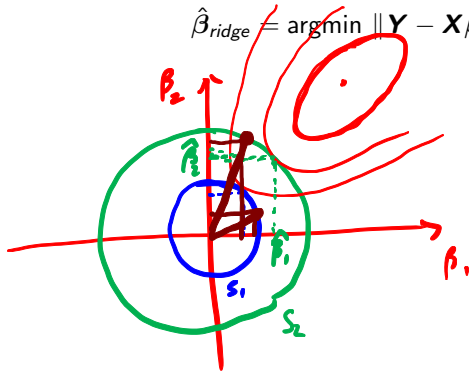
- For any single element: $|\beta_{2j}| \text{ ? } |\beta_{1j}|$, **not sure in general!**
 - In some special cases, we know the relationship for sure, e.g., in class you derived the case when \mathbf{X} is identity matrix.
- For the sum of squares $\|\beta_2\|_2^2 < \|\beta_1\|_2^2$
 - which implies when $p = 1$, $|\beta_2| < |\beta_1|$.

Geometric intuition / alternative view

There is a one-to-one relationship between λ and s in

$$\hat{\beta}_{\text{ridge}} = \operatorname{argmin} \|\mathbf{Y} - \mathbf{X}\beta\|_2^2 + \lambda \|\beta\|_2^2$$

$$\hat{\beta}_{\text{ridge}} = \operatorname{argmin} \|\mathbf{Y} - \mathbf{X}\beta\|_2^2 \quad \text{s.t.} \quad \|\beta\|_2^2 \leq s$$



$$\lambda_1 > \lambda_2$$

$$s_1 < s_2$$

Proof 1: using SVD

In order to show this is true, we need to review something called **singular value decomposition (SVD)**.

- For any real matrix $X \in R^{n \times p}$, $n \geq p$, we can write
- $X = UDV^T$, where
 1. $U \in R^{n \times p}$ and $U^T U = I_p$,
 2. $D \in R^{p \times p}$ and is diagonal,
 3. $V \in R^{p \times p}$ and $V^T = V^{-1}$.

$$X = U D V^T$$

Proof 1: Plug in SVD

$$X = UDV^T \quad U^T U = I, \quad V^T = V^{-T}$$

$$\hat{\beta}_{\text{ridge}} = (X^T X + \lambda I)^T X^T Y$$

$$= (VDU^T \cancel{U} DV^T + \lambda I)^T VDU^T Y$$

$$= (VD^2 V^T + \lambda I)^T VDU^T Y$$

$$= (VD^2 V^T + V(\lambda I)V^T)^T VDU^T Y$$

$$= (V(D^2 + \lambda I)V^T)^T VDU^T Y$$

$$= V(D^2 + \lambda I)^T \cancel{V^T} VDU^T Y$$

$$= V(D^2 + \lambda I)^T D U^T Y \Rightarrow \text{Let } U^T Y \in \mathbb{R}^{p \times 1} = a$$

$$= \sum_{j=1}^p u_j \left(\frac{d_j}{d_j^2 + \lambda} \right) a_j$$

$$\begin{aligned} V^T V &= I \\ \lambda I &= \lambda V V^T \\ &= V(\lambda I)V^T \end{aligned}$$

$$\begin{aligned} (V \tilde{D} V^T)(V \tilde{D}^{-1} V^T) &= V \tilde{D} \tilde{D}^{-1} V^T \\ &= V V^T = I \end{aligned}$$



Proof 1: Plug in SVD

$$\begin{aligned}
 \|\hat{\beta}_{\text{ridge}}\|_2^2 &= \hat{\beta}_{\text{ridge}}^T \hat{\beta}_{\text{ridge}} = Y^T U D (D^2 + \lambda I)^{-1} \cancel{V^T V} (D^2 + \lambda I)^{-1} D U^T Y \\
 &= a^T \text{diag}(\dots) a \\
 &= \sum_{j=1}^p a_j^2 \cdot \frac{d_j^2}{(d_j^2 + \lambda)^2}
 \end{aligned}$$



$$\lambda \uparrow \quad \|\hat{\beta}\|_2^2 \downarrow$$

Proof 2: By definition

Consider the general case with $\lambda_1 < \lambda_2$ (in the previous definition, $\lambda_1 = 0, \lambda_2 = \lambda$)

$$\beta_1 = \operatorname{argmin} \|\mathbf{Y} - \mathbf{X}\beta\|_2^2 + \lambda_1 \|\beta\|_2^2$$

$$\beta_2 = \operatorname{argmin} \|\mathbf{Y} - \mathbf{X}\beta\|_2^2 + \lambda_2 \|\beta\|_2^2$$

Now by definition,

$$\begin{aligned} \|\mathbf{Y} - \mathbf{X}\beta_1\|_2^2 + \lambda_1 \|\beta_1\|_2^2 &\leq \|\mathbf{Y} - \mathbf{X}\beta_2\|_2^2 + \lambda_1 \|\beta_2\|_2^2 \\ \|\mathbf{Y} - \mathbf{X}\beta_2\|_2^2 + \lambda_2 \|\beta_2\|_2^2 &\leq \|\mathbf{Y} - \mathbf{X}\beta_1\|_2^2 + \lambda_2 \|\beta_1\|_2^2 \end{aligned}$$

Adding the two inequalities together,

$$\begin{aligned} \lambda_1 \|\beta_1\|_2^2 + \lambda_2 \|\beta_2\|_2^2 &\leq \lambda_1 \|\beta_2\|_2^2 + \lambda_2 \|\beta_1\|_2^2 \\ (\lambda_1 - \lambda_2) \|\beta_1\|_2^2 &\leq (\lambda_1 - \lambda_2) \|\beta_2\|_2^2 \\ \|\beta_1\|_2^2 &\geq \|\beta_2\|_2^2 \end{aligned}$$

Proof 2: By definition

Now plug in $\|\beta_1\|^2 > \|\beta_2\|^2$ back in to

$$\|\mathbf{Y} - \mathbf{X}\beta_1\|_2^2 + \lambda_1 \|\beta_1\|_2^2 \leq \|\mathbf{Y} - \mathbf{X}\beta_2\|_2^2 + \lambda_1 \|\beta_2\|_2^2$$

$$\begin{aligned} \text{RSS}_1 &\leq \text{RSS}_2 + \underbrace{\lambda_1 \|\beta_1\|_2^2} - \underbrace{\lambda_1 \|\beta_2\|_2^2} \\ &\leq \text{RSS}_2 \end{aligned}$$

Summary so far

So far we have shown that for ridge regression

- When we **increase λ , $\|\beta\|_2^2$ decreases.**
- But we cannot guarantee every element in $|\beta|$ decreases.
- When we **increase λ , RSS increases.**

Now this fact allows us to compare more complicated scenarios...

Ridge: changing scale and λ accordingly

$$\beta_2 = \operatorname{argmin} \|\mathbf{Y} - \mathbf{X}\beta\|_2^2 + \lambda\|\beta\|_2^2$$

$$\beta_5 = \operatorname{argmin} \|\mathbf{Y} - a\mathbf{X}\beta\|_2^2 + a^2\lambda\|\beta\|_2^2$$

By change of variable, if we let $\alpha = a \times \beta$ in the second problem,

$$\min_{\beta} \|\mathbf{Y} - a\mathbf{X}\beta\|_2^2 + a^2\lambda\|\beta\|_2^2 = \min_{\alpha} \|\mathbf{Y} - \mathbf{X}\alpha\|_2^2 + \lambda\|\alpha\|_2^2$$

Notice we know the minimizer of RHS is β_2 , thus $\alpha = \beta_2$, or

$$\beta_5 = \frac{1}{a}\beta_2$$

And RSS is the same.

Ridge: changing scale and λ fixed

$$\beta_2 = \operatorname{argmin} \|\mathbf{Y} - \mathbf{X}\beta\|_2^2 + \lambda\|\beta\|_2^2$$

$$\beta_5 = \operatorname{argmin} \|\mathbf{Y} - a\mathbf{X}\beta\|_2^2 + a^2\lambda\|\beta\|_2^2$$

$$\beta_4 = \operatorname{argmin} \|\mathbf{Y} - a\mathbf{X}\beta\|_2^2 + \lambda\|\beta\|_2^2$$

Consider $a > 1$, using β_5 as a middle step,

- $\|\beta_4\|_2^2 > \|\beta_5\|_2^2 = \frac{1}{a^2}\|\beta_2\|_2^2$
- which implies if $p = 1$, $|\beta_4| > |\beta_5| = \frac{1}{a}|\beta_2|$

($a < 1$ is just the opposite.)

$$\beta_1 = \operatorname{argmin} \|\mathbf{Y} - \mathbf{X}\beta\|_2^2$$

$$\beta_2 = \operatorname{argmin} \|\mathbf{Y} - \mathbf{X}\beta\|_2^2 + \lambda\|\beta\|_2^2$$

$$\beta_3 = \operatorname{argmin} \|\mathbf{Y} - a\mathbf{X}\beta\|_2^2$$

$$\beta_4 = \operatorname{argmin} \|\mathbf{Y} - a\mathbf{X}\beta\|_2^2 + \lambda\|\beta\|_2^2$$

$$\beta_5 = \operatorname{argmin} \|\mathbf{Y} - a\mathbf{X}\beta\|_2^2 + a^2\lambda\|\beta\|_2^2$$

$$\{\beta_0, \beta_6\} = \operatorname{argmin} \|(\mathbf{Y} + 1) - \beta_0 - (\mathbf{X} + 2)\beta\|_2^2$$

$$\{\beta_0, \beta_7\} = \operatorname{argmin} \|(\mathbf{Y} + 1) - \beta_0 - (\mathbf{X} + 2)\beta\|_2^2 + \lambda\|\beta\|_2^2$$

So, in the simple case of $p = 1$, and $a > 1$

$$|\beta_7| = |\beta_2| = a|\beta_5| < a|\beta_4| < a|\beta_3| = |\beta_1| = |\beta_6|$$

If we replace $|\cdot|$ with $\|\cdot\|_2^2$, it holds for $p > 1$ case.

How about RSS?

We already know $RSS_1 = RSS_3 = RSS_6$:

$$\beta_1 = \operatorname{argmin} \|\mathbf{Y} - \mathbf{X}\beta\|_2^2$$

$$\beta_3 = \operatorname{argmin} \|\mathbf{Y} - a\mathbf{X}\beta\|_2^2$$

$$\{\beta_0, \beta_6\} = \operatorname{argmin} \|(\mathbf{Y} + 1) - \beta_0 - (\mathbf{X} + 2)\beta\|_2^2$$

Remember RSS increases as λ increases, we have $RSS_3 < RSS_4 < RSS_5$

$$\beta_3 = \operatorname{argmin} \|\mathbf{Y} - a\mathbf{X}\beta\|_2^2$$

$$\beta_4 = \operatorname{argmin} \|\mathbf{Y} - a\mathbf{X}\beta\|_2^2 + \lambda\|\beta\|_2^2$$

$$\beta_5 = \operatorname{argmin} \|\mathbf{Y} - a\mathbf{X}\beta\|_2^2 + a^2\lambda\|\beta\|_2^2$$

And we also have shown $RSS_2 = RSS_5 = RSS_7$

$$\beta_2 = \operatorname{argmin} \|\mathbf{Y} - \mathbf{X}\beta\|_2^2 + \lambda\|\beta\|_2^2$$

$$\beta_5 = \operatorname{argmin} \|\mathbf{Y} - a\mathbf{X}\beta\|_2^2 + a^2\lambda\|\beta\|_2^2$$

$$\{\beta_0, \beta_7\} = \operatorname{argmin} \|(\mathbf{Y} + 1) - \beta_0 - (\mathbf{X} + 2)\beta\|_2^2 + \lambda\|\beta\|_2^2$$

Questions/exercises for you

- Derive using SVD that RSS can be written as

$$\mathbf{Y}^T \mathbf{Y} + \sum_{j=1}^p \left(\left(\frac{d_j^2}{d_j^2 + \lambda} - 1 \right)^2 - 1 \right) \mathbf{a}_j^2$$

where $\mathbf{a} = \mathbf{U}^T \mathbf{Y}$, and show this is an increasing function of $\lambda > 0$.

- Can you say something similar if we change all ridge regression into Lasso?

Derivation of lasso estimator when $\mathbf{X}^T \mathbf{X} = \mathbf{I}$

To derive the lasso estimator, we first notice that

$$\begin{aligned}
 \hat{\beta}^{(lasso)} &= \operatorname{argmin}_{\beta} \{ (\mathbf{Y} - \mathbf{X}\beta)^T (\mathbf{Y} - \mathbf{X}\beta) + \lambda \sum_{j=1}^p |\beta_j| \} \\
 &= \operatorname{argmin}_{\beta} \{ \mathbf{Y}^T \mathbf{Y} - 2\beta^T \mathbf{X}^T \mathbf{Y} + \beta^T \mathbf{X}^T \mathbf{X} \beta + \lambda \sum_{j=1}^p |\beta_j| \} \\
 &= \operatorname{argmin}_{\beta} \{ \mathbf{Y}^T \mathbf{Y} - 2\beta^T \mathbf{X}^T \mathbf{Y} + \beta^T \beta + \lambda \sum_{j=1}^p |\beta_j| \} \\
 &= \operatorname{argmin}_{\beta} \left\{ \sum_{i=1}^n y_i^2 + \sum_{j=1}^p (\beta_j^2 + \lambda |\beta_j| - 2\beta_j \sum_{i=1}^n x_{ij} y_i) \right\}
 \end{aligned}$$

Derivation of lasso estimator when $\mathbf{X}^T \mathbf{X} = \mathbf{I}$

The above objective function allows us to optimize each element of β separately. Similar to what we have derived in class, we can see the minimum is achieved at

$$\hat{\beta}_j = \max\left\{\sum_{i=1}^n x_{ij}y_i - \frac{\lambda}{2}, 0\right\} \quad \text{if } \hat{\beta}_j > 0$$

$$\hat{\beta}_j = \min\left\{\sum_{i=1}^n x_{ij}y_i + \frac{\lambda}{2}, 0\right\} \quad \text{if } \hat{\beta}_j < 0$$

which gives us the lasso solution

$$\hat{\beta}^{(lasso)} = \begin{cases} \sum_{i=1}^n x_{ij}y_i - \frac{\lambda}{2} & \text{if } \sum_{i=1}^n x_{ij}y_i > \frac{\lambda}{2} \\ \sum_{i=1}^n x_{ij}y_i + \frac{\lambda}{2} & \text{if } \sum_{i=1}^n x_{ij}y_i < -\frac{\lambda}{2} \\ 0 & \text{otherwise} \end{cases}$$