

Ridge regression

Normal regression: $\hat{\theta} = (X^T X)^{-1} X^T Y$

$(p+1) \cdot n$ $n \cdot (p+1)$

Now, $\text{rank}(X) < p+1$, then $\text{rank}(X^T X) < (p+1)$

which means

$$X^T X \Rightarrow \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1(p+1)} \\ 0 & a_{22} & \dots & a_{2(p+1)} \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & \dots & 0 \end{bmatrix}$$

Therefore $(X^T X)^{-1}$ does not exist

Solution: Add a λ on diagonal. to $X^T X$

$$X^T X + \lambda I \Rightarrow \begin{bmatrix} a_{11} & \dots & a_{1(p+1)} \\ \vdots & \ddots & \vdots \\ 0 & \dots & 0 \end{bmatrix} + \lambda \begin{bmatrix} 1 & 0 & \dots & 0 \\ \vdots & \ddots & \ddots & \vdots \\ 0 & \dots & \dots & 1 \end{bmatrix}$$
$$\Rightarrow \begin{bmatrix} a_{11} + \lambda & \dots & a_{1(p+1)} \\ \vdots & \ddots & \vdots \\ 0 & \dots & \lambda \end{bmatrix}$$

$$\text{rank}(X^T X + \lambda I) = p+1$$

Therefore $(X^T X + \lambda I)^{-1}$ does exist.

$$\hat{\theta} = (X^T X + \lambda I)^{-1} X^T Y$$

How to understand ridge?

Going backwards though the proof.

$$SSE = (Y - X\hat{\theta})^T (Y - X\hat{\theta}) + \lambda \theta^T \theta$$
$$\frac{\partial SSE}{\partial \theta} = \frac{\partial [(Y - X\hat{\theta})^T (Y - X\hat{\theta})] + \lambda \partial [\theta^T \theta]}{\partial \theta}$$

\Downarrow

$$\frac{\partial SSE}{\partial \theta} = -2X^T Y + 2[X^T X + \lambda I] \hat{\theta}$$

加上 λI
从下往上
推导

New cost Function

$$SSE = (Y - X\hat{\theta})^T (Y - X\hat{\theta}) + \lambda \theta^T \theta$$
$$= (Y - X\hat{\theta})^T (Y - X\hat{\theta}) + \lambda [\theta_0 \dots \theta_p] \begin{bmatrix} \theta_0 \\ \vdots \\ \theta_p \end{bmatrix}$$
$$= (Y - X\hat{\theta})^T (Y - X\hat{\theta}) + \lambda \sum_{i=0}^p \theta_i^2$$

$$\hat{\theta} = \underset{\hat{\theta}}{\operatorname{argmin}} [(Y - X\hat{\theta})^T (Y - X\hat{\theta}) + \lambda \sum_{i=0}^p \theta_i^2]$$

we want to find a $\hat{\theta}$, such that SSE is minimal.

From Lagrange (拉格朗日优化). The Above could be write as the following

$$\hat{\theta}_{\text{ridge}} = \underset{\hat{\theta}}{\operatorname{argmin}} [(Y - X\hat{\theta})^T (Y - X\hat{\theta})]$$

subject to $\sum_{i=1}^p \theta_i^2 \leq t^2$

where $t^2 \uparrow$ then $\lambda \downarrow$
 $\lambda \uparrow$ then $t^2 \downarrow$

Understand Ridge From graph.

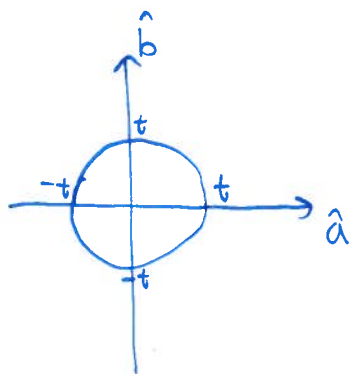
For 2-D Data sets

$$Y = aX + b, \text{ SSE} = \sum (y_i - \hat{a}x_i - \hat{b})^2$$

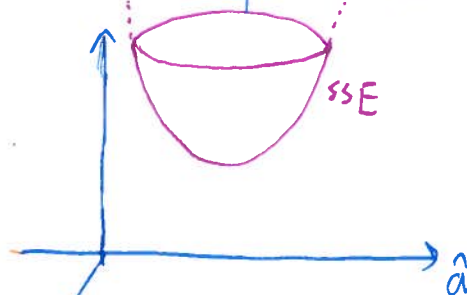
$$\hat{\theta} = \begin{bmatrix} \hat{a} \\ \hat{b} \end{bmatrix} = \underset{(\hat{a}, \hat{b})}{\operatorname{argmin}} \sum (y_i - \hat{a}x_i - \hat{b})^2$$

subject to $\hat{a}^2 + \hat{b}^2 \leq t^2$

$$\hat{a}^2 + \hat{b}^2 = t^2$$



SSE could be proved to be a convex



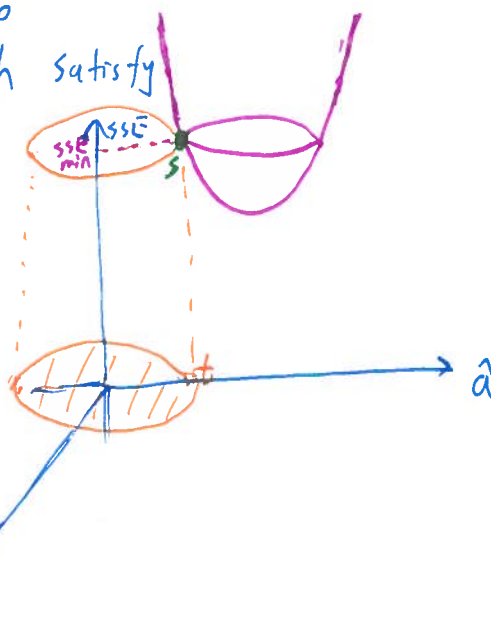
↙ \hat{b}

We are looking for (\hat{a}, \hat{b}) , which satisfy

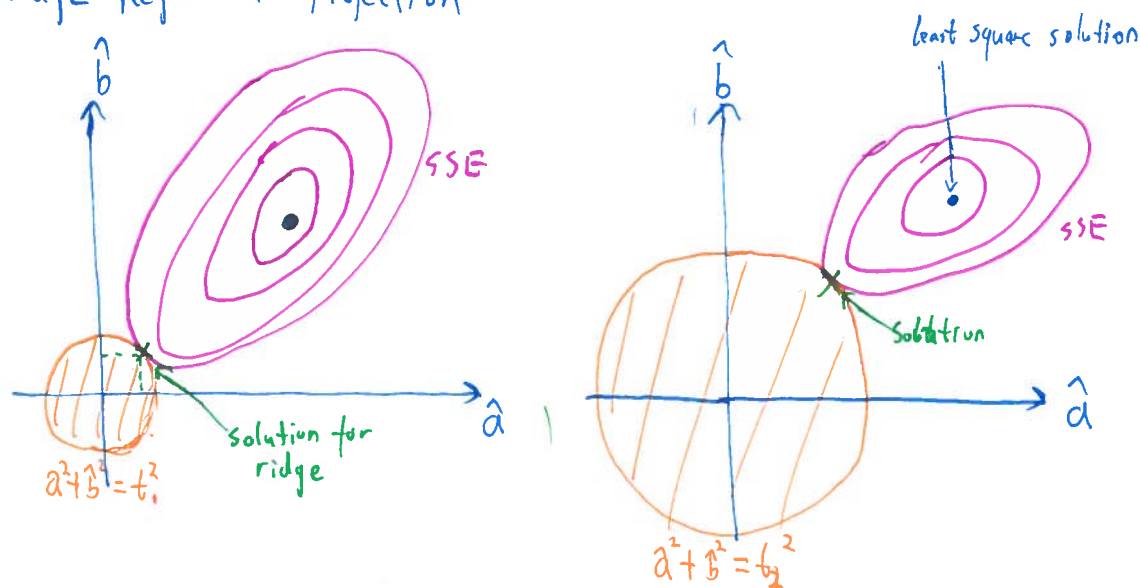
1. $\hat{a}^2 + \hat{b}^2 = t^2$

2. $\text{SSE} = \sum (y_i - \hat{a}x_i - \hat{b})^2$

We could find it in the graph.
point S is the solution



Ridge Regression Projection



$$\hat{\theta}^{\text{ridge}} = \underset{\hat{\theta}}{\operatorname{argmin}} (Y - X\hat{\theta})^T (Y - X\hat{\theta}) + \lambda \sum \theta_i^2$$

$$\hat{\theta}^{\text{ridge}} = \underset{\hat{\theta}}{\operatorname{argmin}} (Y - X\hat{\theta})^T (Y - X\hat{\theta}) \quad \text{subject to } \sum \theta_i^2 \leq t^2$$

Note, when $\lambda \uparrow$, $t^2 \downarrow$
 when $\lambda \downarrow$, $t^2 \uparrow$

when t^2 decrease, λ increase

$$\hat{a} \rightarrow 0, \hat{b} \rightarrow 0$$

when t^2 increase, λ decrease

$[\hat{a}, \hat{b}]$ approach least square solution

Therefore, λ acts like a penalty.

For

$$\hat{\theta} = (X^T X + \lambda I)^{-1} X^T Y$$

when $\lambda \uparrow$, $\hat{\theta}$ approach 0.

when $\lambda \rightarrow 0$, $\hat{\theta}$ approach the true least square solution $(X^T X)^{-1} X^T Y$

Lasso Regression

$$\hat{\theta} = \operatorname{argmin}_{\theta} [(Y - X\theta)^T (Y - X\theta)] + \lambda \sum |\theta_i|$$



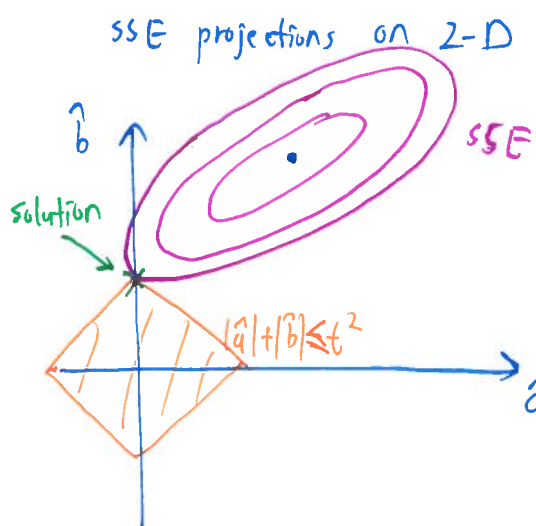
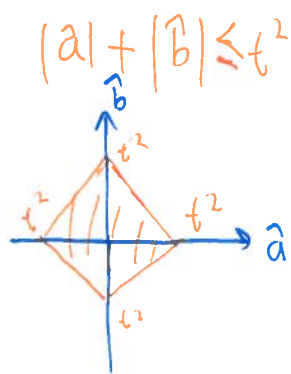
$$\hat{\theta} = \operatorname{argmin}_{\theta} [(Y - X\theta)^T (Y - X\theta)]$$

subject to $\sum |\theta_i| \leq t^2$

when $\lambda \uparrow, t^2 \downarrow$

$\lambda \downarrow, t^2 \uparrow$

For 2-D



Therefore, for Lasso, **SSE** is very likely to touch the pointy point, which implies, some θ_i in the solution are 0.

Solution of Lasso is not discussed here because it is kind of complicated.

Derivation of coordinated descent for Lasso

Lasso cost function [coordinated descent] (交点式证法)

$$\begin{aligned} SSE(\theta) &= \overset{OLS}{SS E}(\theta) + \lambda \|\theta\| \\ &= \frac{1}{2} \sum_{i=1}^n \left[y_i - \sum_{j=0}^p \theta_j x_{ij} \right]^2 + \lambda \sum_{j=0}^p |\theta_j| \end{aligned}$$

where do they come from

OLS:
original-least-square

$$\begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} = \begin{matrix} i: [1, n] \\ j: [0, p] \end{matrix} \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1p} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \vdots & \vdots & x_{ij} & \ddots & \vdots \\ \vdots & x_{n1} & \dots & \dots & x_{np} \end{bmatrix} \times \begin{bmatrix} \theta_0 \\ \vdots \\ \theta_p \end{bmatrix}$$

i : sample j : feature

For i th sample: $y_{pred} = \hat{\theta}_0 + x_{i1} \hat{\theta}_1 + x_{i2} \hat{\theta}_2 + \dots + x_{ip} \hat{\theta}_p$

error² for i th sample: $[y_{true} - y_{pred}]^2$

$$\text{error}^2 \text{ sum} = \sum (y_{true} - y_{pred})^2 = \sum_{i=1}^n \left[y_i - \sum_{j=0}^p \theta_j x_{ij} \right]^2$$

Lasso penalty $\lambda \|\theta\| = \lambda \sum_{j=0}^p |\theta_j|$

$$= \lambda [|\theta_0| + |\theta_1| + |\theta_2| + \dots + |\theta_p|]$$

we want to solve

$$\frac{\partial SSE(\theta)}{\partial \theta} = \frac{\partial \left[\frac{1}{2} \sum_{i=1}^n \left(y_i - \sum_{j=0}^p \theta_j x_{ij} \right)^2 + \lambda \sum_{j=0}^p |\theta_j| \right]}{\partial \theta} = 0$$

OLS part:

$$\begin{aligned}
 \frac{\partial \text{SSE}^{\text{OLS}}(\theta)}{\partial \theta_j} &= \frac{\partial}{\partial \theta_j} \left[\frac{1}{2} \sum_{i=1}^n \left(y_i - \sum_{j=0}^p \theta_j x_{ij} \right)^2 \right] \\
 &= - \sum_{i=1}^n x_{ij} \left[y_i - \sum_{j=0}^p \theta_j x_{ij} \right] \\
 &= - \sum_{i=1}^n x_{ij} \left[y_i - \sum_{k \neq j}^p \theta_k x_{ik} - \theta_j x_{ij} \right] \\
 &= - \sum_{i=1}^n x_{ij} \left[y_i - \sum_{k \neq j}^p \theta_k x_{ik} \right] + \theta_j \sum_{i=1}^n x_{ij}^2 \\
 &\triangleq -p_j + \theta_j z_j
 \end{aligned}$$

where we define p_j and z_j , which the z_j is

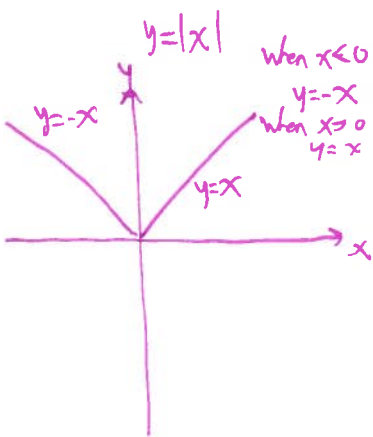
Lasso part:

$$\lambda \sum_{j=0}^p |\theta_j| = \lambda |\theta_j| + \lambda \sum_{k \neq j}^n |\theta_k|$$

$$\frac{\partial [\lambda \sum |\theta|]}{\partial \theta_j} = \frac{\partial [\lambda |\theta_j|]}{\partial \theta_j} + \frac{\partial [\lambda \sum_{k \neq j}^n |\theta_k|]}{\partial \theta_j}$$

$$= \frac{\partial [\lambda |\theta_j|]}{\partial \theta_j} + 0$$

$$= \begin{cases} -\lambda & \text{if } \theta_j < 0 \\ [-\lambda, \lambda] & \text{if } \theta_j = 0 \\ \lambda & \text{if } \theta_j > 0 \end{cases}$$



Putting OLS and Lasso together

$$\frac{\partial SSE(\theta)}{\partial \theta_j} = -p_j + \theta_j z_j + \frac{\partial \lambda \|\theta_j\|}{\partial \theta_j} = 0$$

$$0 = \begin{cases} -p_j + \theta_j z_j - \lambda & \text{if } \theta_j < 0 \\ [-p_j - \lambda, -p_j + \lambda] & \text{if } \theta_j = 0 \\ -p_j + \theta_j z_j + \lambda & \text{if } \theta_j > 0 \end{cases}$$

note: $\theta \in [-p_j - \lambda, -p_j + \lambda]$

$$-\lambda \leq p_j \leq \lambda.$$

final solution

$$\begin{cases} \theta_j = \frac{p_j + \lambda}{z_j} & \text{for } p_j < -\lambda \\ \theta_j = 0 & \text{for } -\lambda \leq p_j \leq \lambda \\ \theta_j = \frac{p_j - \lambda}{z_j} & \text{for } p_j > \lambda \end{cases}$$

$$\text{where } p_j = \sum_{i=1}^n x_{ij} \left[y_i - \sum_{k \neq j}^p \theta_k x_{ik} \right]$$

$$z_j = \sum_{i=1}^n x_{ij}^2$$

How to do Lasso using python

Do it using step-wise way.

def solution(ϕ_j, λ):

$$\text{Initial } \theta = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} \text{ or } \theta = \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix}$$

$$\theta_j = \begin{cases} \phi_j + \lambda & , \phi_j < -\lambda \\ 0 & , -\lambda \leq \phi_j \leq \lambda \\ \phi_j - \lambda & , \phi_j > \lambda \end{cases}$$

for $\lambda = 0.01$ in $[1, 50000]$:

for each col j in X :

$$X_{-j} = X[:, j] = \begin{bmatrix} x_{1j} \\ x_{2j} \\ \vdots \\ x_{nj} \end{bmatrix}$$

$$Y_{\text{pred}} = X \cdot \theta = \begin{bmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_n \end{bmatrix}$$

$$\phi_j = \sum_{i=1}^n x_{ij} \left[y_i - \sum_{k \neq j} \theta_k x_{ik} \right]$$

$$= [x_{1j} \ x_{2j} \ \dots \ x_{nj}] \left\{ \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} - \begin{bmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_n \end{bmatrix} \right\}$$

(when $k=j$)

$$= [x_{1j} \ \dots \ x_{nj}] \left\{ \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} - \begin{bmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_n \end{bmatrix} + \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_j \\ \vdots \\ 0 \end{bmatrix} \right\}$$

$$= [x_{1j} \ \dots \ x_{nj}] \left\{ \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} - \begin{bmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_n \end{bmatrix} + \theta_j \begin{bmatrix} x_{1j} \\ \vdots \\ x_{nj} \end{bmatrix} \right\}$$

$$= X_{-j}^T \cdot (y - \hat{y} + \theta[j] \cdot X_{-j})$$

$$\theta[j] = \text{Solution}(\phi_j, \lambda)$$