

Shrinkage Methods (Ridge and Lasso)

Petchara Pattarakijwanich

Introduction to Data Science, 28 October 2022

Goal of this week

- Subset Selection
 - Best Subset Selection
 - Forward/Backward Stepwise Subset Selection
- Shrinkage Methods
 - Ridge Regression
 - The Lasso

About the Final Project

Approaches

- Answer specific questions from data, using data science techniques
- Implement, invent, or test some data science techniques or algorithms

Data Source

- Your own research
- Publicly available datasets
<https://geekflare.com/open-datasets-for-data-science/>
<https://www.dataquest.io/blog/free-datasets-for-projects/>
- Simulated data (?)

Presentation (some time around final week or a bit later)

Model Selection

ให้แบบจำลองที่ดีที่สุด

$$\begin{aligned} \boxed{\text{AIC}} &= \frac{1}{\sigma^2} \left[\boxed{\text{MSE}} + \boxed{\frac{2d\sigma^2}{n}} \right] \\ \boxed{\text{BIC}} &= \frac{1}{\sigma^2} \left[\boxed{\text{MSE}} + \boxed{\frac{d\sigma^2 \ln n}{n}} \right] \end{aligned}$$

↑
ค่าคงที่ = ไม่สนใจ

↑
พจน์ที่ 1

↑
พจน์ที่ 2

Subset Selection

ตัวอย่าง Input $\{x_1, x_2, x_3\}$ output y
[สมมติว่า y ขึ้นกับ x_1 อย่างเดียว ไม่ขึ้นกับ x_2, x_3]

Linear Regression ในการ subset ของ Input

ๆ ได้มาดังนี้

- 0 ตัวแปร

$$y = \beta_0$$

- 1 ตัวแปร

$$y = \beta_0 + \beta_1 x_1, \quad y = \beta_0 + \beta_2 x_2, \quad y = \beta_0 + \beta_3 x_3$$

- 2 ตัวแปร

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2, \quad \dots, \dots$$

- 3 ตัวแปร

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3$$

การลดตัวแปร 7 ตัวแปร

Best Subset Selection

รู้กันไว้ว่า $y = \beta_0 + \beta_1 x_1$ ดีสุด

- ฝึกหัด 8 โมเดล \Rightarrow คำนวณ AIC, BIC \Rightarrow ฝึกหัด = ดีสุด
(ใช้ test data)

[ถ้า มี p ตัวแปร จำนวนโมเดลที่ต้องฝึก = 2^p]

\Rightarrow ฝึกหมดแล้วกัน ถ้า p มากพอ

[Brute force]

Forward Step-wise Subset Selection

0 parameter $y = \beta_0$

1 parameter $y \propto x_1$, $y \propto x_2$, $y \propto x_3$

- มี 3 ตัว เข้ามา เปรียบเทียบ

- เลือก $y \propto x_1$ ดี

2 parameter $y \propto (x_1, x_2)$, $y \propto (x_1, x_3)$

- มี 2 ตัว (แล้ว $y \propto (x_2, x_3)$)

- เลือก $y \propto (x_1, x_2)$ ดี

3 parameter $y \propto (x_1, x_2, x_3)$

\Rightarrow เมื่อ 4 ตัว มา เปรียบเทียบ AIC, BIC \Rightarrow เลือก ดี

Forward Step-wise Subset Selection

มี p ตัวแปร

0 parameter : 1 โมเดล

1 parameter : p โมเดล \Rightarrow ได้ข้อสรุปว่า ตัวแปร

2 parameters : $p-1$ โมเดล \Rightarrow ได้ข้อสรุปว่า ตัวแปร

3 parameters : $p-2$ โมเดล

\vdots

$p-1$ parameters : 2 โมเดล

p parameters : 1 โมเดล

\uparrow Linear Regression model

จำนวนโมเดลทั้งหมด

$$= 1 + p + (p-1) + (p-2) + \dots + 2 + 1$$

$$= 1 + \frac{1}{2} p(p+1)$$

$$\propto p^2$$

Backward Step-wise Subset Selection

အဲဒါ Forward လေးက ပဲ အဲဒါပဲ \Rightarrow အဲဒါပဲ 1

Ridge Regression

↑ Regularization technique

LR : Minimize $RSS = \sum (y - y_{\text{model}})^2$

$$= \sum (y - \beta_0 - \beta_1 x_1 - \beta_2 x_2 - \dots)^2$$

$$RSS = \sum (y - \beta_0 - \sum_i \beta_i x_i)^2$$

Ridge Regression :

Minimize

$$\sum (y - \beta_0 - \sum_i \beta_i x_i)^2 + \lambda \sum_i \beta_i^2$$

RSS

penalty

regularization = $\beta_1^2 + \beta_2^2 + \dots$
= two-norm

Ridge Regression

Ridge Regression ;

Minimize

$$\underbrace{\sum (y - \beta_0 - \sum_i \beta_i x_i)^2}_{\text{RSS}} + \underbrace{\lambda \sum_i \beta_i^2}_{\text{penalty}}$$

= $\beta_1^2 + \beta_2^2 + \dots$
= two-norm

$\Rightarrow \lambda$ is the tuning parameter

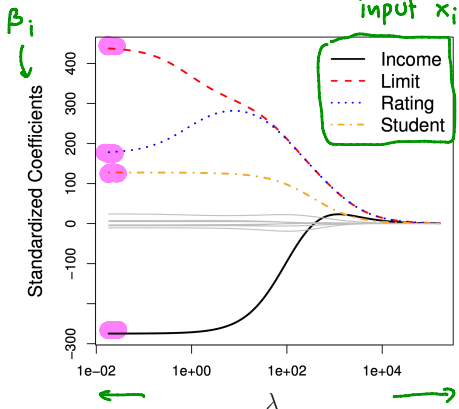
\Rightarrow Minimize also

- ↗ Minimize RSS (RSS ថែវ)
- ↘ β_i ធំតិច ($\sum_i \beta_i^2$ ថែវ)

\Rightarrow Shrinkage method
(Regularization)

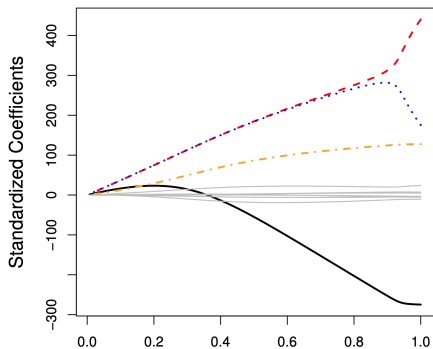
Ridge Regression

$$\|\beta\|_2 = \text{two-norm} = \sum_i \beta_i^2$$



$\lambda \rightarrow 0$
is just LR

$\lambda \rightarrow \infty$
 $\beta_i = 0$
 $y = \beta_0 = \bar{y}$



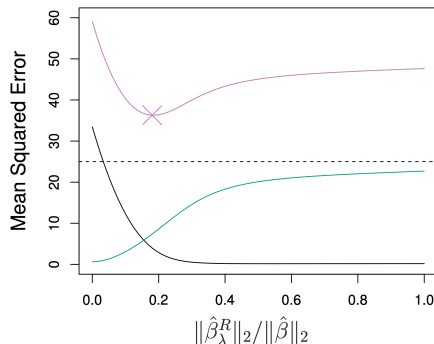
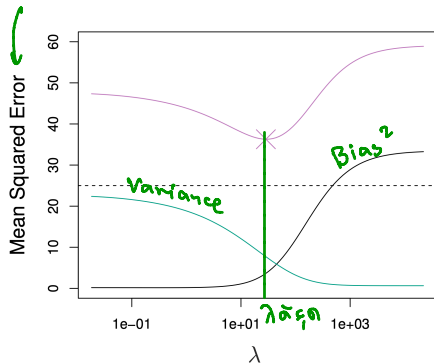
$\lambda \rightarrow \infty$
 $\beta_i = 0$

$\|\hat{\beta}_\lambda^R\|_2 / \|\hat{\beta}\|_2 = \frac{(\sum \beta_i^2)_{\text{Ridge}}}{(\sum \beta_i^2)_{\text{LR}}}$

$\lambda \rightarrow 0$
LR

Ridge Regression

vs test data



as λ increases \Rightarrow noisy test data

λ increases \Rightarrow model stiff

- Bias increases
- Variance decreases

\Rightarrow Bias Variance trade off

The Lasso

LR : Minimize $RSS = \sum (y - y_{\text{model}})^2$

$$= \sum (y - \beta_0 - \beta_1 x_1 - \beta_2 x_2 - \dots)^2$$

$$RSS = \sum (y - \beta_0 - \sum_i \beta_i x_i)^2$$

Lasso :

Minimize

$$\sum (y - \beta_0 - \sum_i \beta_i x_i)^2 + \lambda \sum_i |\beta_i|$$

regularization $= |\beta_1| + |\beta_2| + \dots$
 $= \text{one-norm}$

RSS

penalty

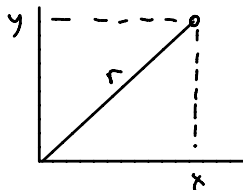
The Lasso

Two-norm $\|x\|_2$ vs

สองนอร์ม 2 นอร์ม

(Euclidean distance)

$$r^2 = x^2 + y^2 = \|x\|_2^2$$

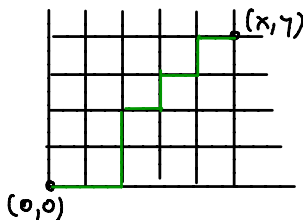


One-norm $\|x\|_1$

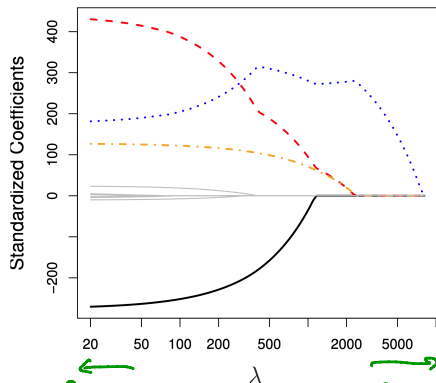
สองนอร์ม 1 นอร์ม

(Taxi-driver distance)

$$r = |x| + |y| = \|x\|_1$$

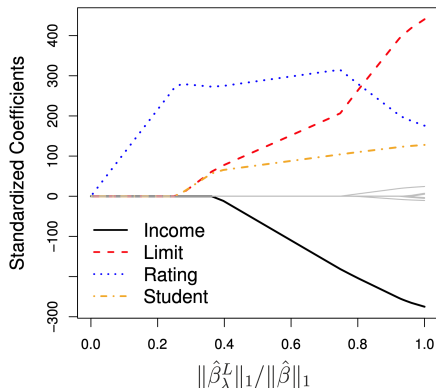


The Lasso



←
 $\lambda \rightarrow 0$
 OLS LR

→
 $\lambda \rightarrow \infty$
 $\beta_i = 0$
 $y = \beta_0 = \bar{y}$



Parameter Selection : เมื่อ λ มากพอแล้ว β_i หมดทั้ง = 0 ทีเดียว

The Lasso

LR : Minimize
$$RSS = \sum (y - \beta_0 - \sum_i \beta_i x_i)^2$$

Ridge : Minimize
$$\sum (y - \beta_0 - \sum_i \beta_i x_i)^2 + \lambda \sum_i \beta_i^2$$

Lasso : Minimize
$$\sum (y - \beta_0 - \sum_i \beta_i x_i)^2 + \lambda \sum_i |\beta_i|$$

LR သာ (သုတေသန) x_i များကို အသုံးပြုခြင်း \Rightarrow β_i များကို အသုံးပြုခြင်း x_i

\Rightarrow Ridge & Lasso သာ x_i များကို အသုံးပြုခြင်း $\Rightarrow \beta_i$ များကို အသုံးပြုခြင်း ပြုသော ပြုသော

\Rightarrow သို့မဟုတ် normalize သာ x_i \nearrow β_i များကို $[0, 1]$
 \searrow β_i mean = 1, SD = 0

Ridge vs Lasso

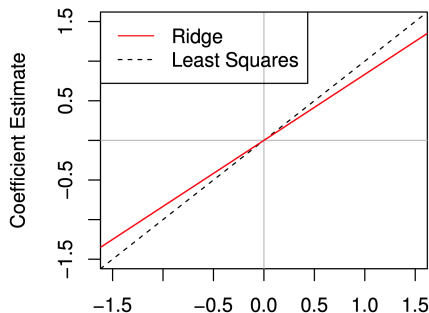
ວິທີໃຫຍ່ຕົ້ນ ໑ ເລືອກ λ ທີ່ເໝາະສົມແນວໃດ ?

- ການຮອບກັບ test data (ນັກ cross validation)

ວິທີໃຫຍ່ ອັດຕາໄລ່

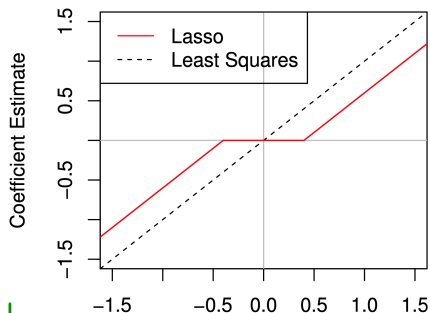
- Ridge ອັດຕາໄລ່ input ທຸກໆຕົວແທນ output ($\beta_i \neq 0$)
 - Lasso ອັດຕາໄລ່ input ທຸກໆຕົວແທນ output (β_i ລາຍງາຍ = 0 ຮ່າງກາຍ ແກງຍາວ)
- ນັກ parameter selection ໂດຍອັດຕາໄລ່

Ridge vs Lasso



Ridge

$$\beta_{i,\text{ridge}} \approx \frac{\beta_{i,\text{LR}}}{C(\lambda)}$$



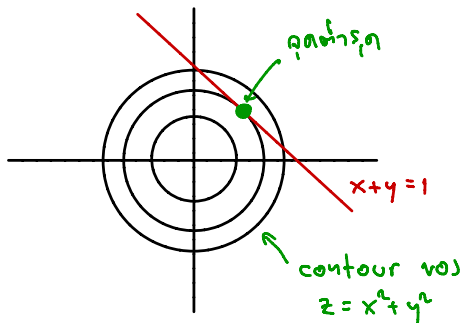
Lasso

$$|\beta_{i,\text{lasso}}| \approx \begin{cases} |\beta_{i,\text{LR}}| - C(\lambda) \\ 0 \end{cases}$$

Ridge vs Lasso

Lagrange Multiplier

Ex. หา (x, y) ที่ทำให้ $\underbrace{z = x^2 + y^2}_{\text{objective function}}$ มีค่าต่ำสุด $\underbrace{x + y = 1}_{\text{constraint}}$



Ridge vs Lasso

Ex. ឃើញ (x, y) អំពី $z = x^2 + y^2$ តាម $x + y = 1$

$$\underbrace{x + y = 1}_{\phi(x, y) = 0}$$
$$x + y - 1 = 0$$

ឃើញ (x, y) អំពី $z = x^2 + y^2$ តាម $\phi = 0$

(Constrained optimization problem)

⇓ Lagrange Multiplier

ឃើញ (x, y, λ) អំពី $z + \lambda \phi$ តាម

⇓
(Unconstrained optimization problem)

Ridge vs Lasso

Ridge

$$(y - \beta_0 - \sum_i \beta_i x_i)^2 + \lambda \sum_i \beta_i^2 \iff$$

$$\begin{aligned} &\text{minimize } (y - \beta_0 - \sum_i \beta_i x_i)^2 \\ &\text{s.t. } \sum_i \beta_i^2 = d \lambda \end{aligned}$$

Lasso

$$(y - \beta_0 - \sum_i \beta_i x_i)^2 + \lambda \sum_i |\beta_i| \iff$$

$$\begin{aligned} &\text{minimize } (y - \beta_0 - \sum_i \beta_i x_i)^2 \\ &\text{s.t. } \sum_i |\beta_i| = d \lambda \end{aligned}$$

Ridge vs Lasso

