

Week 9: predict CLTV

Suppose we want to predict $CLTV_i$ (transaction in one year)

using recency R_i

R_i : days since last transaction.

R_i captures a customer is active or not.

$E[CLTV_i | R_i]$ is the best predictor.

Q: how to estimate it?

Sample analogue

$$\hat{E}[CLTV | R_i = r]$$

= sample average of CLTV with $R_i = r$

Quiz: how would you calculate $\hat{E}[CLTV | R_i = r]$?

e.g. $\hat{E}[CLTV | R_i = 7]$

subsample last transact on 2022/10/24

Y_i = total transaction between 2022/10/31 ~ 2023/10/31.

We can estimate $E[CLTV | R_i]$ by OLS regression:

sample avg. of $R_i = n$ $\mathbb{1}$ is its indicator function.

$$CLTV_i = \beta_1 \mathbb{1}\{R_i = 1\} + \beta_2 \mathbb{1}\{R_i = 2\} + \dots + \beta_n \mathbb{1}\{R_i = n\}$$

① convenient package

② easy-to-communicate

③ inference

Ex: Is vaccine effective?

Covid: infected or not?

Vaccine: Vaccinated or not?

Sup vaccine is randomly given

We can compare

$$\overline{\text{Covid}}_{\text{Vacc}=1} - \overline{\text{Covid}}_{\text{Vacc}=0} \dots \textcircled{1}$$

You can also:

$$\text{Covid} = \beta_0 + \beta_1 \text{Vaccine} + \varepsilon$$

and $\textcircled{1}$ will be your $\hat{\beta}_1$

$$\text{Covid} = \gamma_0 + \gamma_1 \text{Vaccine} + \gamma_2 \text{Male} + \gamma_3 \text{Vaccine} \cdot \text{Male}$$

$\gamma_1 + \gamma_3$ is the effect of vaccine on male.

Very often, estimators in econometrics are equivalent to an OLS w/

certain specification OLS is a work horse.

We know:

$$\text{CLTV} = \beta_1 \mathbb{1}\{R_i=1\} + \beta_2 \mathbb{1}\{R_i=2\} + \dots + \beta_p \mathbb{1}\{R_i=p\}$$

estimate the conditional mean.

The equation has 90 parameters if $n = 1800$ (small) When p large, we call it "high dimensional"

$$\frac{n}{p} = \frac{1800}{90} = 20 \Rightarrow \text{Use 20 obs. to estimate one mean} \Rightarrow \text{noisy}$$

The ratio $\frac{n}{p}$ determines whether the model is high dim.

$\frac{n}{p}$: > 1 moderate
 < 1 very high
 $>> 1$ low dim

Alternative model:

$$\begin{aligned} \text{CLTV}_i = & r_1 \mathbb{1}\{1 \leq R_i \leq 7\} \\ & + r_2 \mathbb{1}\{8 \leq R_i \leq 15\} \\ & + \dots + \\ & + r_{14} \mathbb{1}\{85 \leq R_i \leq 100\} \end{aligned}$$

↓

$$\text{Now, } \frac{n}{p} = \frac{100}{14} > 100$$

easier to estimate, but biased $\therefore \neq$ conditional mean.

Alternatively,

$$\begin{aligned} \text{CLTV} = & b_0 \mathbb{1}\{1 \leq R_i \leq 20\} \\ & + b_1 \mathbb{1}\{21 \leq R_i \leq 60\} \\ & + b_2 \mathbb{1}\{61 \leq R_i \leq 90\} \end{aligned}$$

easier to estimate, but more biased.

Daily, $p=90$, $\frac{n}{p} \approx 20$, non-biased

Weekly $p=14$, $\frac{n}{p} > 100$, medium-biased

Monthly $p=3$, $\frac{n}{p} > 500$, highly-biased.

Facing a trade-off problem:

bias-variance tradeoff.

bias: whether the model is similar to best predictor.

variance: whether the model is easy to estimate.

can be chosen by cross-validation.

"All models are wrong, but some are useful"

Penalized regression

last class "smoothing" exploiting the continuity of $f(r) = E[CLTV | R_i = r]$

Penalization

Penalize complicated models (model that have many non-zero coefficients)

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_p X_{ip}$$

if only a few β_j is not zero \Rightarrow simple model.

Penalized regression.

$$\min_{\beta \in \mathbb{R}^p} \frac{1}{2n} \sum (y_i - x_i^T \beta)^2 + \lambda \|\beta\|_q$$

penalization.

where $\lambda > 0$ is the penalty term.

$$\|\beta\|_q = \sqrt[q]{\sum \|\beta_j\|^q}$$

Remark:

λ is a "framing parameter" specified by the user.

"There are many penalized regression methods $q=1$. LASSO (Least Absolute Selection and Shrinkage Operator)

Why $q=1$?

① convexity (for $q \geq 1$)

if $q < 1$, objective function isn't convex.

eg. $q = \frac{1}{2}$. $p=1$ $\|\beta\| = \sqrt{\beta}$, not convex.

② $g=1$ has sparsity 稀疏 Γ α

$$\hat{\beta}_{\text{Lasso}}: \begin{pmatrix} \hat{\beta}_1 \\ \vdots \\ \hat{\beta}_p \end{pmatrix} \text{ most } \hat{\beta}_j = 0$$

Advantage of sparsity:

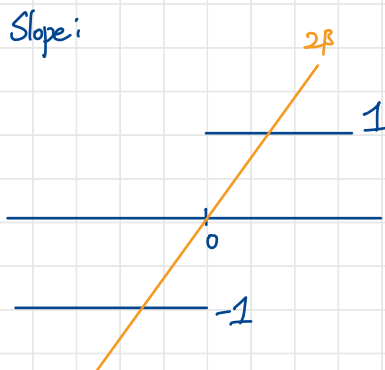
① easy to interpret.

② Can be used for variable selection.
(only choose X_j if $\hat{\beta}_j \neq 0$)

Penalty ($p=1$)

$$\|\beta\|_1 = |\beta|$$

β^2



Marginal penalty:

move β away from 0

marginal penalty = 1 $\Rightarrow g=1$

" $\approx 0 \Rightarrow g=2$

LASSO regression likes to stay in 0, different from Ridge.

$$p=2, \beta = \begin{pmatrix} \beta_1 \\ \beta_2 \end{pmatrix}$$

What is a circle?

$$g=1 \quad \|\beta\|_1 = |\beta_1| + |\beta_2|$$

$$\{x \mid \|x\|_1 = 1\}$$

$$g=2 \quad \|\beta\|_2 = \sqrt{\beta_1^2 + \beta_2^2}$$

$$\{x \mid \|x\|_2 = 1\}$$

LASSO

$$\min_{\beta \in \mathbb{R}^p} \sum_{i=1}^n (y_i - x_i^T \beta)^2 + \lambda \|\beta\|_1, \quad \|\beta\|_1 = |\beta_1| + |\beta_2| + \dots + |\beta_p|$$

is differentiable everywhere.

$$\hat{\beta}_{\text{Lasso}} = ?$$

In general, LASSO does not have analytic solution.

Can be written as func. of data.

$$(\hat{\beta}_{\text{OLS}} = (X^T X)^{-1} X^T Y)$$

⇒ has to use numerical optimization.
(versions of gradient descent)

Special Case:

$X_{1i}, X_{2i}, \dots, X_{pi}$ are dummy variables

$$\sum_{i=1}^p X_{1i} = 1$$

e.g.

$$X_{1i} : \mathbb{1}\{\text{Female}\}$$

$$X_{2i} : \mathbb{1}\{\text{Male}\}$$

⇒ (As group membership indicators)

Objective function

$$\sum (y_i - x_i^T \beta)^2 + \lambda \|\beta\|_1$$

$$= \underbrace{\sum_{X_{1i}=1} (y_i - \beta_1)^2}_{L_1} + \underbrace{\sum_{X_{2i}=1} (y_i - \beta_2)^2}_{L_2} + \dots + \sum_{X_{pi}=1} (y_i - \beta_p)^2 + \lambda (|\beta_1| + |\beta_2| + |\beta_3| + \dots + |\beta_p|)$$

$$L_1 = L_1$$

$$L_2 = L_2$$

Define

$$L_j(\beta_j) = \sum_{X_{ij}=1} (y_i - \beta_j)^2 + \lambda |\beta_j|$$

$$\text{LASSO Obj.} = \sum_{j=1}^p L_j(\beta_j)$$

$$\min_{\beta \in \mathbb{R}^p} \sum (y_i - x_i \beta)^2 + \lambda \|\beta\|_1$$

$$\Leftrightarrow \min_{\beta_1} L_1(\beta_1) \\ \min_{\beta_2} L_2(\beta_2) \\ \vdots \\ \min_{\beta_p} L_p(\beta_p)$$

Define

$$n_j = \sum_{i=1}^n \mathbb{1}_{\{x_{ij}=1\}}$$

= # observations

$$x_{ij}=1$$

$$\bar{y}_j = \frac{1}{n_j} \sum_{i: x_{ij}=1} y_i$$

= avg. of y among $x_{ij}=1$

$$L_j(\beta_j) = \sum_{i: x_{ij}=1} (y_i - \beta_j)^2 + \lambda |\beta_j|$$

$$\begin{aligned} \sum (y_i - \beta_j)^2 &= \sum (y_i^2 - 2\beta_j y_i + \beta_j^2) \\ &= \sum y_i^2 - 2\beta_j \sum y_i + n\beta_j^2 \end{aligned}$$

$$= \beta_j^2 - 2n\bar{y}_j \beta_j + c$$

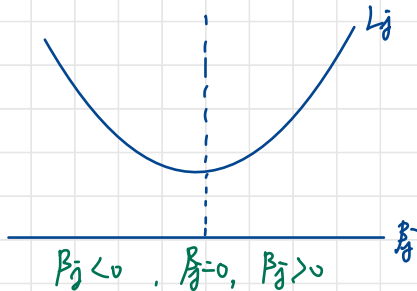
$$\hat{\beta}_j^{OLS} = \bar{y}_j \text{ since}$$

$$n\beta_j^2 - 2n\bar{y}_j \beta_j + c$$

$$= n(\beta_j - \bar{y}_j)^2 + c$$

$$\hat{\beta}_j^{OLS} = \hat{\beta}_j^{OLS} = \bar{y}_j$$

$$L_j(\beta_j) = n(\beta_j - \bar{y}_j)^2 + \lambda |\beta_j|$$



$$\text{At } \beta_j = 0$$

$$L_j(\beta_j)|_{\beta_j=0} = n\bar{y}_j^2$$

$$\text{At } \beta_j > 0$$

$$\frac{dL_j}{d\beta_j} = 2n(\beta_j - \bar{y}_j) + \lambda = 0$$

$$2n\beta_j = 2n\bar{y}_j - \lambda$$

$$\beta_j^* = \bar{y}_j - \frac{\lambda}{2n}, \text{ if } \bar{y}_j - \frac{\lambda}{2n} > 0$$

$$\text{If } \beta_j < 0$$

$$L_j = n(\beta_j - \bar{y}_j)^2 - \lambda \beta_j \quad \frac{dL_j}{d\beta_j} = 0$$

$$\Rightarrow \beta_j = \begin{cases} \bar{y}_j + \frac{1}{2n}\lambda, & \text{if } \bar{y}_j + \frac{\lambda}{2n} < 0 \\ 0 & \end{cases}$$

$$\Rightarrow \hat{\beta}_j = \begin{cases} \bar{y}_j - \frac{1}{2n}\lambda, & \bar{y}_j - \frac{1}{2n}\lambda > 0 \\ 0 & \text{o.v.} \\ \bar{y}_j + \frac{1}{2n}\lambda, & \bar{y}_j + \frac{1}{2n}\lambda < 0 \end{cases}$$

Compare OLS & LASSO

$$\hat{\beta}^{OLS} > \frac{1}{2n}\lambda \Rightarrow \hat{\beta}^{LASSO} = \hat{\beta}^{OLS} - \frac{1}{2n}\lambda$$

$$\hat{\beta}^{OLS} < \frac{1}{2n}\lambda \Rightarrow \hat{\beta}^{LASSO} = \hat{\beta}^{OLS} + \frac{1}{2n}\lambda$$

$$\frac{1}{2n}\lambda \leq \hat{\beta}^{OLS} \leq \frac{1}{2n}\lambda \Rightarrow \hat{\beta}^{LASSO} = 0$$

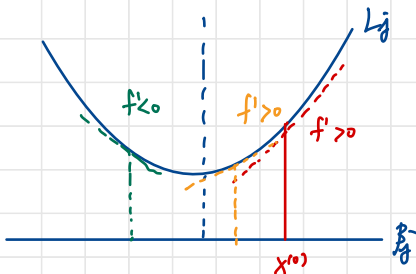
Selection: $\hat{\beta} \rightarrow 0$

only keep coefficients with large $|\hat{\beta}_j^{OLS}|$

Shrinkage: shrink non-zero coefficients toward 0.

Solving LASSO with numerical optimization.

LASSO can be solved by versions of gradient descent.



Start with initial point $x^{(0)}$

Iterate

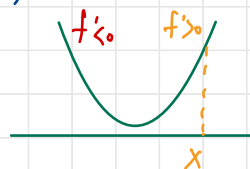
$$x^{(t+1)} = x^{(t)} - \eta^{(t)} df(x^{(t)})$$

$\eta^{(t)}$ step size.

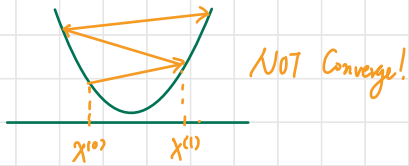
E.g.

$$f(x) = 2(x-1)^2 + 5$$

$$f'(x) = 4(x-1)$$



η too large $\eta_t \rightarrow \infty$



η too small



How to choose $\eta^{(t)}$?

Newton's method.

$$\eta^{(t)} = (\nabla^2 f)^{-1} = (f''(x^{(t)}))^{-1}, \text{ when } p=1$$

→ More conservative (smaller step) when slope is changing fast.

Challenges.

non-convexity: converges to local minima.