

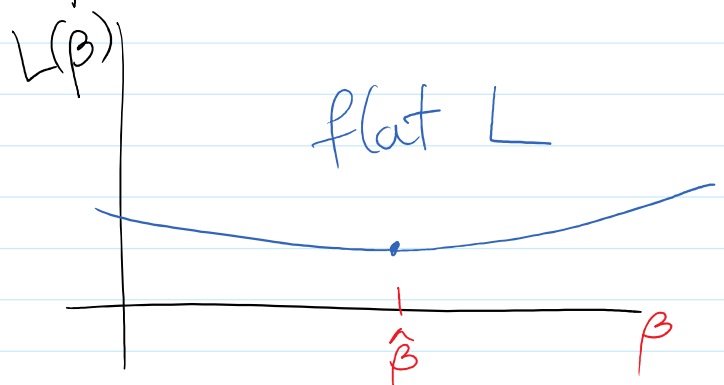
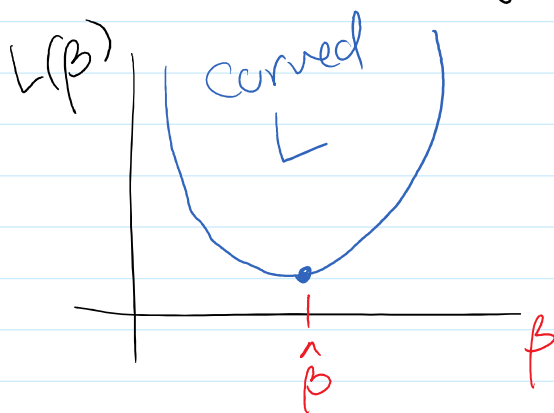
Variable Selection

- All procedures we've looked at allow for some set of P features (regr. or class.)
- May want to: select a "best" set of features
 - ① prediction accuracy
(bias \uparrow and var \downarrow by $P \downarrow$) \otimes
 - ② interpretation
 - ③ model may be ill-conditioned
($P \gg N$)

note: apply to class. and regr

Review of LS Regr.

- $\hat{\beta}$ comes from solving $(X^T X)\beta = X^T Y$
so stability of $\hat{\beta}$ depends on $K(X^T X)$



$K(X^T X)$ small

$K(X^T X)$ large

Simple illustration:

Consider

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots$$

Q: what happens if $X_1 \approx X_2$.

our model:

$$\begin{aligned} Y &\approx \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots \\ &= \beta_0 + (\beta_1 + \beta_2) X_1 + \dots \end{aligned}$$

if $\hat{\beta}_1 = 5$ and $\hat{\beta}_2 = 10$

then just as good of a fit w/ $\hat{\beta}_1 = 7, \hat{\beta}_2 = 8$

b/c $\hat{\beta}_1 + \hat{\beta}_2 = 15$ still.

Generally any $\hat{\beta}_1 + \hat{\beta}_2 = 15$ gives same fit

es. $\hat{\beta}_1 = -1000, \hat{\beta}_2 = 1015$

So when $K(X^T X)$ we tend to have problems

when β s run off to $\pm \infty$.

Variable Selection: choose X_1 , not exclude X_2
from our fit. $\text{Fix } K(X^T X)$

Idea: try a bunch of models w/ different subsets of variables and choose best model.

Way 1: penalize training metric by P

[Classic
very fast]

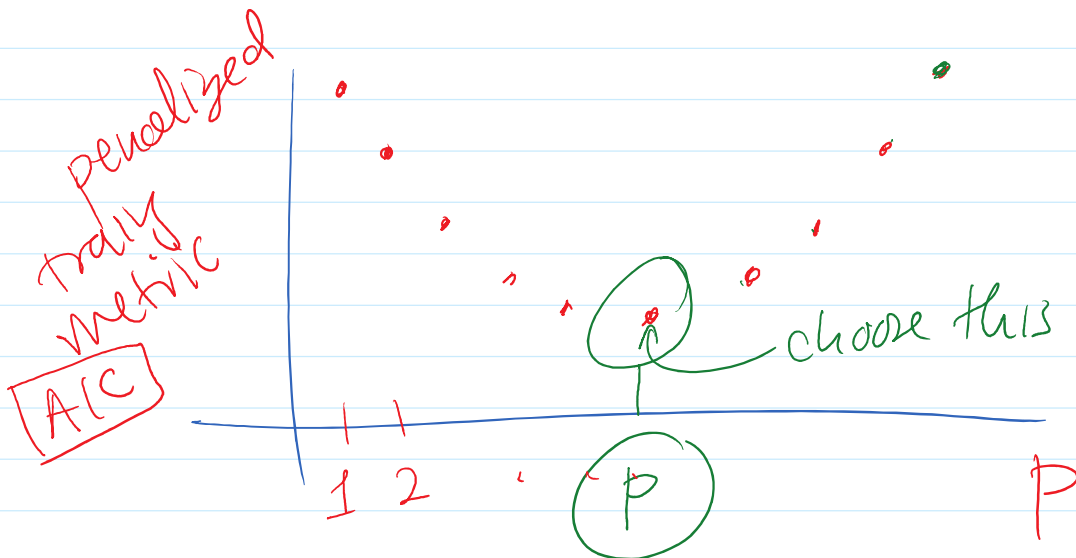
and use this to choose P

Ex. forward stepwise regression

→ add "best" variables 1 at a time

Ex. backwards stepwise

→ start w/ all, remove "worst" var.
1 at a time



Way 2: X-validation

(train/validate/test)

[modern
slaves]

Note: typical testing all possible subsets of vars is not feasible.

We not want to test all possible subsets

Remember: how well does our model building procedure work?

① testing all possible subsets is a very flexible building procedure
(high flex - low bias high var)

② using a constrained/simpler search gives a less flexible model building procedure
(low flex - higher bias lower var)

If our vars are standardized and \pm get $\hat{\beta}$ s order by size (abs)

$$|\hat{\beta}|_{(1)}, |\hat{\beta}|_{(2)}, \dots, |\hat{\beta}|_{(p)}$$

\uparrow \uparrow

$| \hat{\beta}(1) |, | \hat{\beta}(2) |, \dots, | \hat{\beta}(p) |$
 \uparrow largest $\quad \quad \quad \uparrow$ smallest
 $\quad \quad \quad \perp$
 $\quad \quad \quad$ threshold

Do var selection by incl. only covs w/

$$| \hat{\beta}(i) | \geq t$$

Hard-Thresholding

$$\hat{\beta}_i^{(HS)} = \begin{cases} \hat{\beta}_i & \text{if } | \hat{\beta}_i | \geq t \\ 0 & \text{else} \end{cases}$$

Ridge Regression:

Q: Can we deal w/ ill conditioning in a "continuous" way?

Recall: For OLS (Ordinary Least-Squares)

$$L(\beta) = \| y - X\beta \|^2 \text{ and}$$

$$\hat{\beta}^{(OLS)} = \arg \min_{\beta} L(\beta)$$

Ridge Regression: slightly alter OLS
to penalize if p is too large.

Interp 1
of Ridge

$$\hat{\beta}^{(Ridge)} = \arg \min_{\beta} L(\beta) + \lambda \|\beta\|^2$$

size of β

λ = how much we penalize
big β s

$\lambda = 0$ no penalty

$$\hat{\beta}^{(Ridge)} = \hat{\beta}^{(OLS)}$$

$$\lambda \rightarrow \infty, \hat{\beta} \rightarrow 0$$

$\lambda \geq 0$
penalty parameter

typically we don't
penalize β_0

$$\|\beta\|^2 = \sum_{j=1}^p \beta_j^2$$

(secret: choose λ by X -validation)

Ridge Interp 2

Recall Calc III (Lagrange Multipliers)

$$(*) \min_x f(x) \quad \text{s.t.} \quad g(x) \leq t \quad \left[\begin{array}{l} \text{constrained} \\ \text{optim} \end{array} \right]$$

Solved w/ Lagrange multipliers

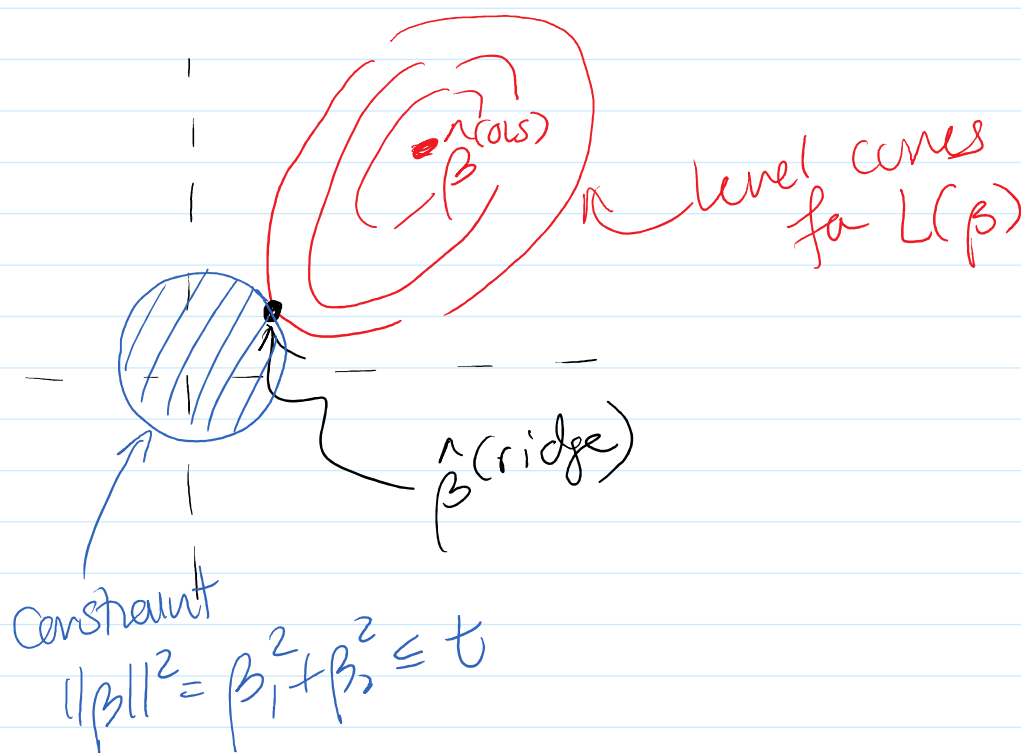
$$\mathcal{L}(\lambda, x) = \underbrace{f(x)} + \underbrace{\lambda(g(x) - t)}$$

ad (x) equiv to
 $\min \mathcal{L}(\lambda, x)$ for some λ
 corresp. to t

Ridge equiv. formulation

$$\hat{\beta}^{(\text{ridge})} = \underset{\beta}{\operatorname{argmin}} L(\beta) \text{ st. } \|\beta\|^2 \leq t$$

corresp. to λ



Interp 3

Note that $L(\beta) = \|y - X\beta\|^2$

Interp 3 | Note that $L(\beta) = \|y - X\beta\|^2$
and $\|\beta\|^2$ are both quadratic,
and differentiable.

For OLS: this means closed form soln

$$\hat{\beta}^{(OLS)} = (X^T X)^{-1} X^T y \quad \left[(X^T X) \beta = X^T y \right]$$

Ridge: also has a closed form soln

$$\hat{\beta}^{(ridge)} = (X^T X + \lambda I)^{-1} X^T y \quad \left[(X^T X + \lambda I) \beta = X^T y \right]$$

replaced $X^T X$ w/ $X^T X + \lambda I$

For OLS: conditioning depended on $K(X^T X)$

Ridge: conditioning depends on $K(X^T X + \lambda I)$

$X^T X$ may not be well cond.

but

$X^T X + \lambda I$ may be better conditioned

Decompose $X = UDV^T$. (SVD) OVS

$$\rightarrow X^T X = V D U^T U D V^T = V D^2 V^T$$

$$\rightarrow (X^T X)^{-1} = V D^{-2} V^T$$

$$\rightarrow (X^T X)^{-1} X^T = V D^{-2} V^T V D U^T = V D^{-1} U^T$$

$$\rightarrow \hat{Y} = X \hat{\beta} = X (X^T X)^{-1} X^T Y = U D V^T V D^{-1} U^T Y$$

$$= U U^T Y$$

$$= \sum_j U_j U_j^T Y$$

als. of U
basis of $\text{Col}(X)$
sum of Y proj. onto U_j

Re-do for Ridge

$$\hat{Y}_{\text{ridge}} = X \hat{\beta}_{\text{ridge}}$$

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

$$(AB)^{-1} = B^{-1} A^{-1}$$

$$= X (V D^2 V^T + \lambda I)^{-1} X^T Y$$

$$= U D V^T (V D^2 V^T + \lambda I)^{-1} V D U^T Y$$

$$= U D (V^T (V D^2 V^T + \lambda I) V)^{-1} D U^T Y$$

$$= U D (D^2 + \lambda I)^{-1} D U^T Y$$

diagonal

$$D = \text{diag}(\sigma_i^2)$$

$$\begin{aligned} & \text{diagonal} \\ & \text{diag}(\sigma_i^2 + \lambda) \\ & \text{inverse} \text{diag}\left(\frac{1}{\sigma_i^2 + \lambda}\right) \end{aligned}$$

$$= U \text{diag}\left(\frac{\sigma_i^2}{\sigma_i^2 + \lambda}\right) U^T Y$$

$$= \sum_j \left(\frac{\sigma_j^2}{\lambda + \sigma_j^2} \right) u_j u_j^T Y$$

weighted sum of proj. of
Y onto u_j

$$\frac{\sigma_i^2}{\lambda + \sigma_i^2} \leq 1 \quad (\text{Shrinkage})$$

Shrink directions w/ lower σ_i more.

Var sel

