

Lecture 12: Variable Selection

So far procedures have a fixed set of vars.

May want to select a "best" set

Why?

① prediction accuracy

② interpretation

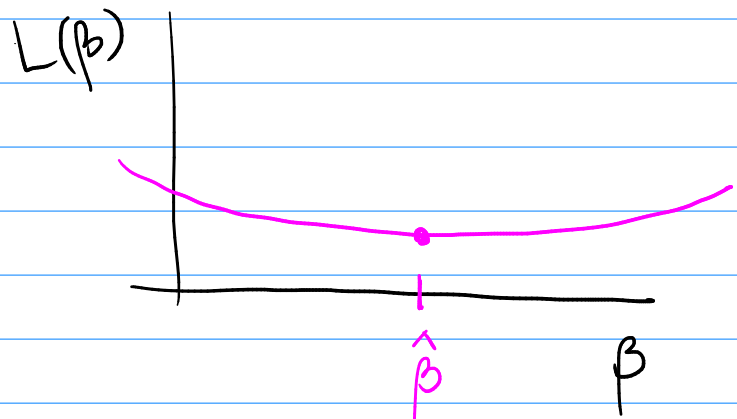
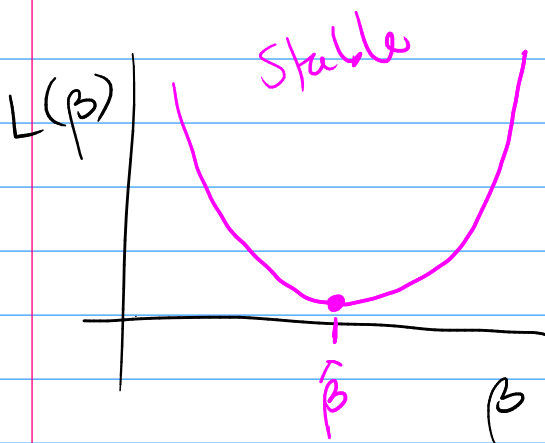
③ model may be ill-conditioned
(OLS regression: $P > N$)

Back to OLS

Recall that $\hat{\beta}$ is obtained by solving

$$X^T X \beta = X^T y$$

the stability of $\hat{\beta}$ depends on inverting $X^T X$



Condition Number

For a linear system $Az = b$

the stability of the solution depends on the condition Number of A denoted $K(A)$
↳ how sensitive is solution to changes in A or b .

Fact: $K(A) = \frac{\sigma_{\max}(A)}{\sigma_{\min}(A)}$ ← largest singular value of A
↑ smallest sing. val.

Large $K(A)$ = very sensitive linear system

Small $K(A)$ = not so sensitive

$K(A) = \infty \Leftrightarrow A$ isn't invertible.

Why do we care?

Want to solve $\underbrace{(X^T X)}_A \underbrace{\beta}_z = \underbrace{X^T y}_b$

The stability of $\hat{\beta} = (X^T X)^+ X^T y$ depends on $K(X^T X)$

If $K(X^T X)$ is large then the regression is ill-conditioned — very sensitive to X and Y

Several causes:

Ex. one var is (approx.) a LC of others
(vars. highly correlated)

Ex. $P > N$ (so $K(X^T X) = \infty$)

e.g. X measures $P = 20K$ genes for
 $N = 30$ patients

How to deal w/ this?

① Variable selection

② Shrinkage

Goal of ① is to pick some subset of "important" variables to use.

→ how do we define?

→ how do we pick?

Two approaches:

① use same individual metric for each var and choose vars w/ best metric.

e.g. calc. p-val. for each var.

and use set of vars w/ smallest p-values.

potential problem: performance of one var. may depend on others.

② calc. same metric on groups of vars - choose group w/ best metric.

potential problem: w/ P vars I have 2^P possible subsets

Careful not to look at training metric

$RSS_{\text{train}} \downarrow$ as $P \uparrow$

Solns: ① X-val. (comp. expensive)

② penalized training metric (classic)

Ex. Adjusted R^2 :

$$R^2_{\text{adj.}} = 1 - \frac{N-1}{N-p-1} (1 - R^2)$$

as $p \uparrow$ eventually $R^2_{\text{adj.}} \downarrow$

Ex. Mallows's C_p

$$C_p = \frac{1}{N} (RSS_{\text{train}} + \underbrace{2p\hat{\sigma}^2}_{\text{penalty}})$$

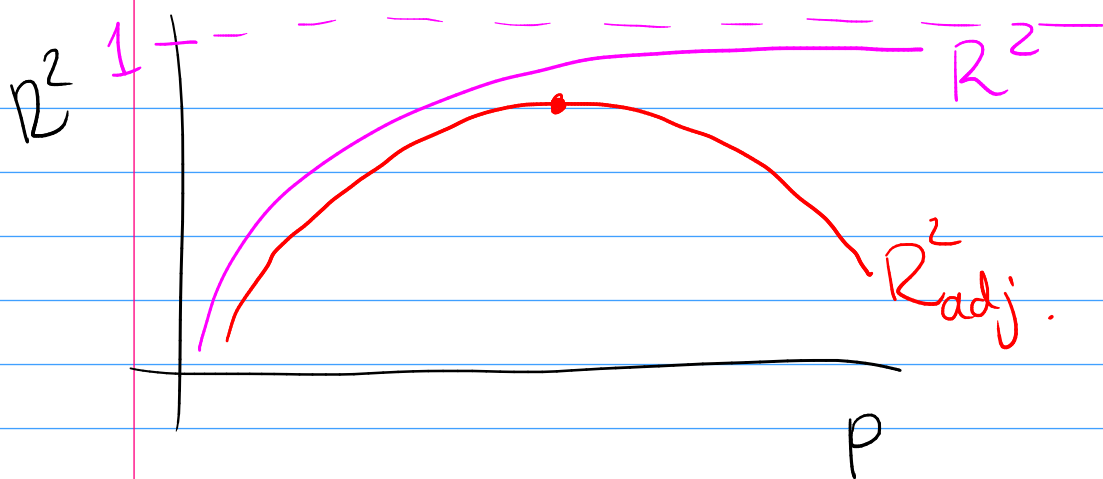
Ex. AIC

$$AIC = \frac{1}{N\hat{\sigma}^2} (RSS_{\text{train}} + 2p\hat{\sigma}^2)$$

Ex. BIC

$$BIC = \frac{1}{N} (RSS_{\text{train}} + \log(N)p\hat{\sigma}^2)$$



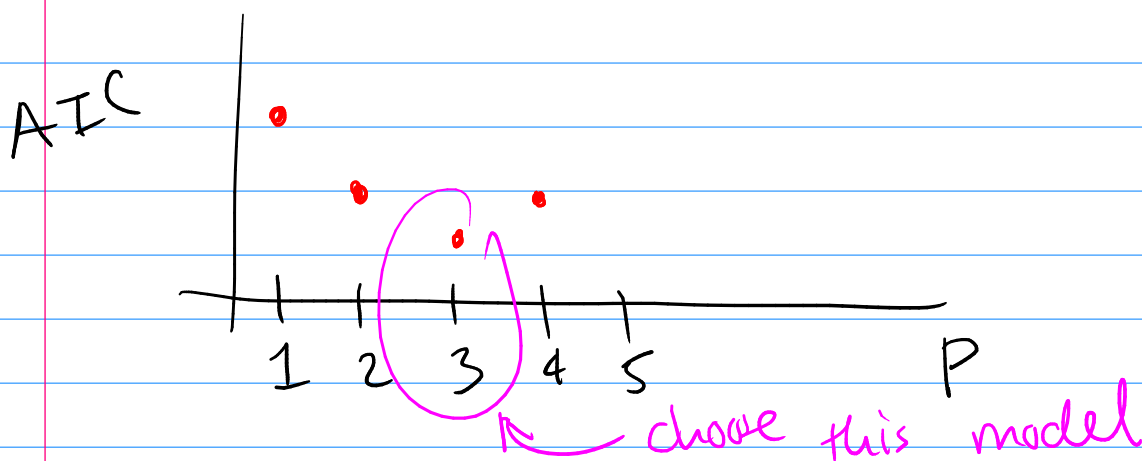


Problem is that potentially have 2^p models to consider

Use a greedy approach:

① Forward Selection:

- ① start w/ model w/ just intercept
- ② add variable that improves metric the most (inc. R^2_{adj} or dec. AIC)
- ③ repeat step ② until my metric gets worse



(2) Backwards Selection

- start w/ model w/ all vars
- removes them one-at-a-time.

Can I deal w/ ill-conditioning in a more continuous way?

Ridge Regression

For OLS we minimize

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

i.e. $\hat{\beta} = \underset{\beta}{\operatorname{argmin}} L(\beta)$

If some of my vars are highly correlated then the assoc. coef. tend to blow-up ($\pm\infty$)

Ex. $Y \approx \hat{\beta}_0 + \hat{\beta}_1 X_1 + \hat{\beta}_2 X_2$ say $\hat{\beta}_1 = 5$, $\hat{\beta}_2 = 7$

if $X_1 \approx X_2$ then model is

$$Y \approx \hat{\beta}_0 + (\hat{\beta}_1 + \hat{\beta}_2) X_1$$

equiv. as good 'is' $\hat{\beta}_1 = 5000$, $\hat{\beta}_2 = -4988$

Ridge regression penalizes optimization to avoid large $\hat{\beta}$ s.

$$\hat{\beta}^{(\text{ridge})} = \underset{\beta}{\operatorname{argmin}} L(\beta) + \lambda \|\beta\|^2$$

$\lambda \geq 0$

By adding an $\lambda \|\beta\|^2$ if elements of β get too large this penalty becomes large.

If $\lambda = 0$ we get OLS $\hat{\beta}$.

as $\lambda \rightarrow \infty$ we get $\hat{\beta}^{(\text{ridge})} \rightarrow 0$

(Typically exclude β_0 from penalty)

Often we standardize vars before ridge.

Also typically choose λ via X-val.