So far procedures have had a fixed set of vars.

May want to select "best" set:

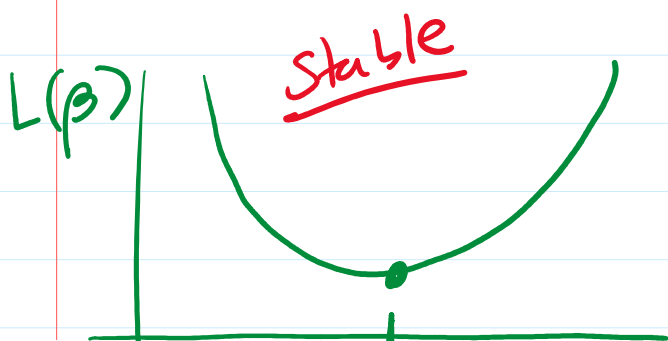Why? ① prediction accuracy : bias and variance

② interpretation

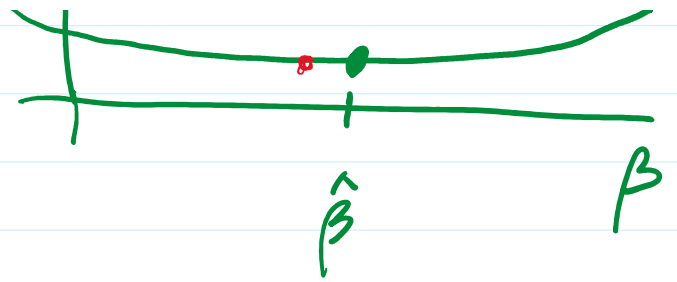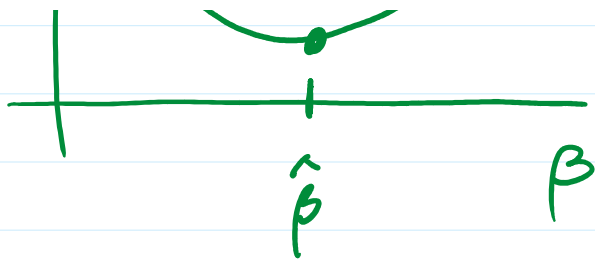## Back to OLS

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

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

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

$L(\beta)$     Stable

$L(\beta)$     unstable

---

## Condition Number

For a linear system $Az = b$ the stability of soln deps on the **condition number** of $A$:

$$K(A).$$

Basically: measures how sens. soln. is to small perturbations in $A$ or $b$.

Fact: $K(A) = \dfrac{\sigma_{max}(A)}{\sigma_{min}(A)}$ ← largest sing. val. of $A$

smallest ↗

Large $K(A)$ = very sensitive (ill-conditioned)

Small $K(A)$ = insensitive (well-cond.)

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

---

Why do I care?

Want to solve $(X^T X)\beta = X^T y$.

$\underbrace{X^T X}_{A} \quad \underbrace{\beta}_{z} \quad \underbrace{X^T y}_{b}$

then the stability of $\hat{\beta}$ depends on $K(X^T X)$.

---

Why do we get a large $K(X^T X)$?

① If $P < N$ but one var. is (approx)
   a lin. comb. of others

② If $P > N$ then $K(X^T X) = \infty$

   e.g. $X$ meas. $P = 20,000$ genes
   among $N = 30$ patients

How can we deal w/ this?

　① Variable selection

　② Shrinkage

---

Goal of ① is to pick some subset of important vars to use.

Approach #1 calc. some importance metric for each var and then only keep those w/ best value.

e.g. calc. p-val. for each $\hat{\beta}_i$ and only retain those w/ low p-vals.

problem: perf. of one var may depend on others.

- - - - - - - - - - - -

Approach #2: calc. metric for groups of

Approach #2: calc. metric for groups of vars and choose group w/ best val.

problem: w/ p vars I have $2^P$ possible subsets to try.

---

Possible metrics for approach #2

Know: don't look at traing metric alone as $P \uparrow$ the traing metric $\downarrow$.
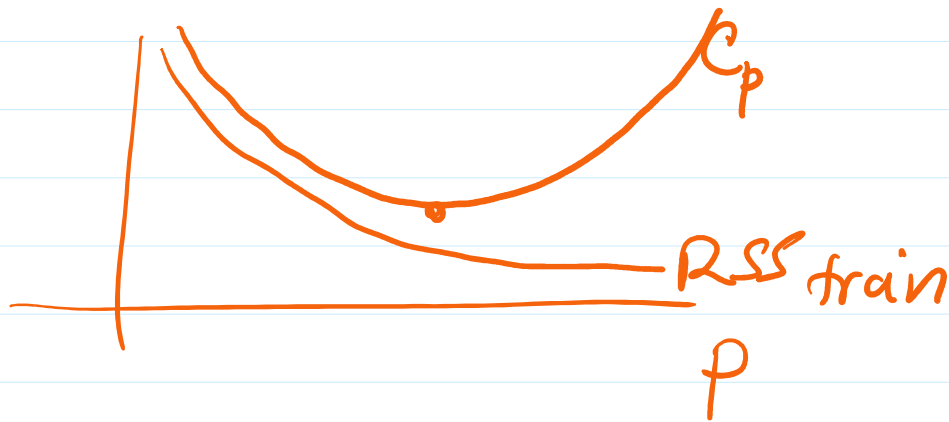
Soln! ① use some val/testy data or cross-val efe.

② penalized traing metric (classic)

$$RSS_{train} = \sum_n (y_n - \hat{y}_n)^2$$

Ex. Mallow's $C_p$

## Ex. Mallow's $C_p$

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



## Ex. AIC $= \frac{1}{N\hat{\sigma}^2}\left(RSS_{train} + 2P\hat{\sigma}^2\right)$

## Ex. BIC $= \frac{1}{N}\left(RSS_{train} + \log(N)P\hat{\sigma}^2\right)$

— — — — — — — — —

## Ex. Adjusted $R^2$

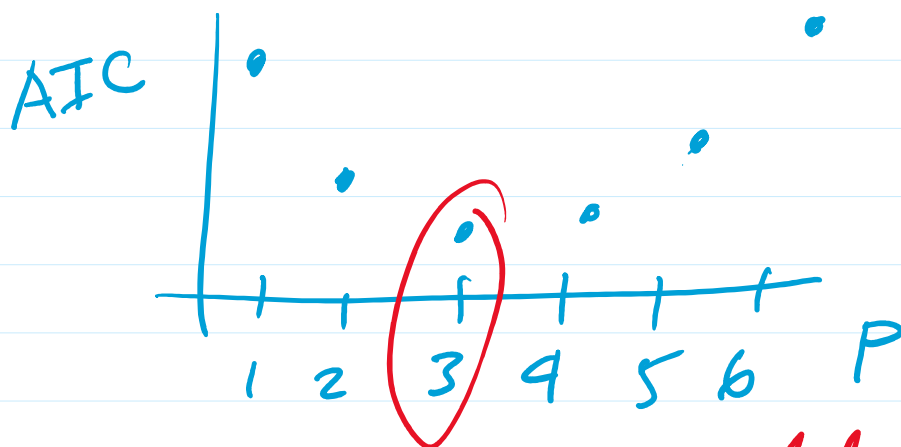$$R^2_{adj} = 1 - \frac{N-1}{N-P-1}\left(1 - R^2\right)$$

Problem: have $2^P$ models to check

Problem: have $2^r$ models to check

Use a ~~greedy~~ approach

## Forward Selection

(i) start w/ model w/ just intercept

(ii) add var. to model that improves metric most
  (dec. AIC   or inc. Adj. $R^2$)

(iii) repeat (ii) until my metric stops improving.



Choose this model

Can I deal w/ ill-conditioned problems in a continuous way? _Shrinkage_.

## Ridge Regression

For OLS we minimize

$$L(\beta) = RSS(\beta) = \|y - X\beta\|_2^2$$

and we let

$$\hat{\beta} = \underset{\beta}{\text{argmin}}\, L(\beta).$$

If my $X^TX$ is ill conditioned (some of my vars are highly cor) and then my elements of $\hat{\beta}$ tend to get really large ($\rightarrow \pm \infty$)

_Ex._ $Y \approx \hat{\beta}_0 + \hat{\beta}_1 X_1 + \hat{\beta}_2 X_2$

but $X_1 \approx X_2$          $\quad\}$     $\hat{\beta}_1 = 5$

but $X_1 \approx X_2$

Say $\hat{\beta}_1 = 5$
$\hat{\beta}_2 = 7$

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

$$\approx \hat{\beta}_0 + \underbrace{(\hat{\beta}_1 + \hat{\beta}_2)}_{12} X_1$$

basically as good as

$$\hat{\beta}_1 = 5000 \qquad \hat{\beta}_2 = -4988$$