

Lecture 13: Variable Selection

All the procedures so far have worked w/ a fixed set of features

May want to select the "best" set of features.
Why?

- ① prediction accuracy
 - ② interpretation
 - ③ model may be ill-conditioned
($P \gg N$)
-

Back to LS OLS regression:

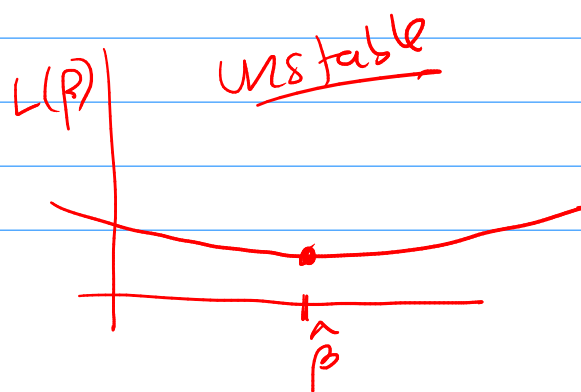
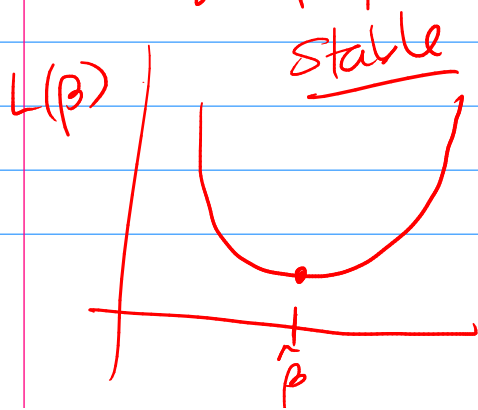
recall that $\hat{\beta}$ comes from solving "normal" eqns

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

to get $\hat{\beta}$ we multiply by $(X^T X)^{-1}$ to get

$$\hat{\beta} = (X^T X)^{-1} X^T y$$

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



If $X^T X$ isn't invertible, $\hat{\beta}$ not unique.

Simple illustration

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

What happens if $X_1 \approx X_2$ (highly correlated)

then

$$Y \approx \beta_0 + \beta_1 X_1 + \beta_2 X_1 + \dots$$

$$\approx \beta_0 + (\beta_1 + \beta_2) X_1 + \dots$$

If $\beta_1 = 5, \beta_2 = 7$ \uparrow (sum = 12)

then basically as good to have $\beta_1 = 0, \beta_2 = 12$
(sum = 12)

or $\beta_1 = -100, \beta_2 = 112$

What tends to happen is that $\beta_1, \beta_2 \rightarrow \pm\infty$

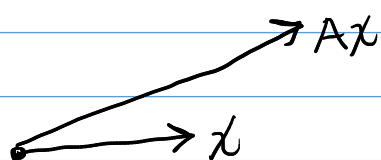
The stability depends on inverting $X^T X$

i.e. the condition number of $X^T X$.

Condition Number:

For A ($n \times n$ matrix) let

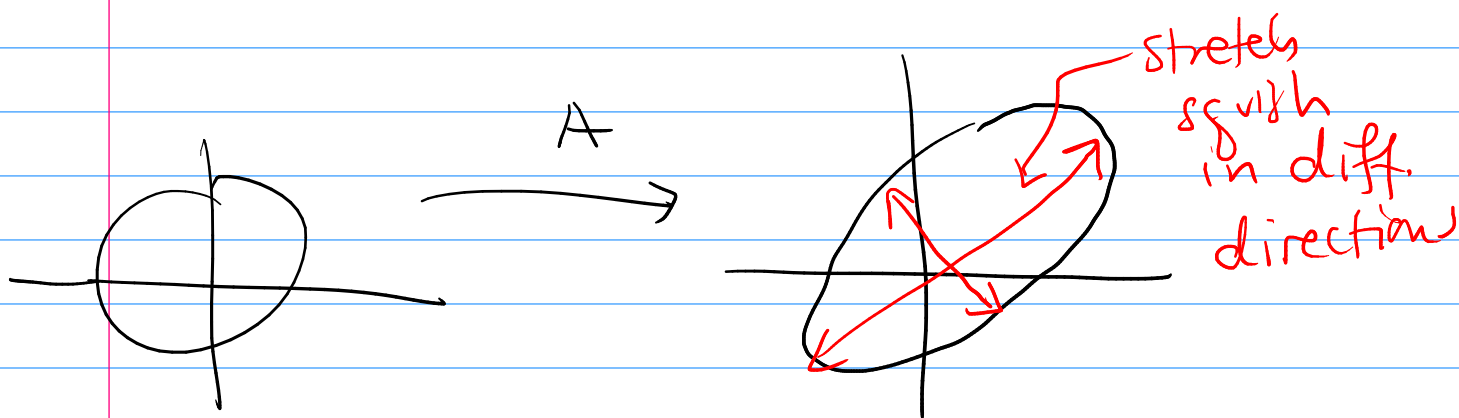
$$M = \max_x \frac{\|Ax\|}{\|x\|} = \max_{x: \|x\|=1} \|Ax\|$$



= max amt. A stretches a unit vector x

$$m = \min_x \frac{\|Ax\|}{\|x\|} = \min_{\|x\|=1} \|Ax\|$$

= min. amt. A stretches a unit vector.



The condition number $K(A) = \frac{M}{m}$

Imagine solving $Az = b$

$$A(z + \delta z) = (b + \delta b)$$

↖ perturb b by δb

my soln is perturbed to $z + \delta z$

So I have $Az = b$ and $A\delta z = \delta b$

notice: $\textcircled{1} \|b\| = \|Az\| \leq M \|z\|$ (by defn)

$$\| \delta b \| = \| A \delta z \| \geq m \| \delta z \| \quad \left[\frac{\|Az\|}{\|z\|} \leq M \right]$$

$$\left[\frac{\|A \delta z\|}{\|\delta z\|} \geq m \right]$$

rearrange $\textcircled{1} \quad \frac{M}{\|b\|} \geq \frac{1}{\|z\|}$

$$\textcircled{2} \quad \frac{\|\delta b\|}{m} \geq \|\delta z\|$$

multiply

$$\underbrace{\frac{M}{m}}_{K(A)} \underbrace{\frac{\|\delta b\|}{\|b\|}}_{\substack{\downarrow \\ \text{rel. size} \\ \text{of } \delta b}} \geq \underbrace{\frac{\|\delta z\|}{\|z\|}}_{\substack{\downarrow \\ \text{rel. size} \\ \text{of } \delta z}}$$

If I perturb b by δb , rel. change in z may be up to $K(A)$ times rel. change in b .

If $K(A)$ is small then system is stable,
b/c rel. large change in b produce
small changes in z

and vice-versa.

Notice that if A isn't invertible then

$$m = 0$$

so $K(A) = \infty$

Fact: $K(A) = \frac{\sigma_{\max}}{\sigma_{\min}}$ \leftarrow larger sing. val.,
 \leftarrow smallest sing. val.

pf. $A = UDV^T$

$$M = \max_{\|x\|=1} \|Ax\| = \max_{\|x\|=1} \|UDV^T x\|$$

Claim: multiplying by orthog. mtrx
doesn't change length of vector

$\|a\| = \sqrt{a^T a} \rightarrow$

$$\|Qx\| = \sqrt{(Qx)^T (Qx)} = \sqrt{x^T Q^T Q x} = \sqrt{x^T x} = \|x\|$$

$$= \max_{\|y\|=1} \|Uy\|$$

$$= \max_{\|y\|=1} \|Dy\| \quad D = \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_n \end{bmatrix}$$

$\sigma_1 \geq \sigma_2 \geq \sigma_3 \geq \dots$

$$= \max_{\|y\|=1} \sqrt{(\sigma_1 y_1)^2 + (\sigma_2 y_2)^2 + \dots + (\sigma_n y_n)^2}$$

choose $y = (1, 0, 0, \dots, 0)$

$$= \sqrt{\sigma_1^2} = \sigma_1 = \sigma_{\max}$$

Similarly $m = \sigma_{\min}$

Why do we care?

For regression solve $(X^T X)\beta = X^T y$
and so the stability of system depends on
 $K(X^T X)$.

We can get a large $K(X^T X)$ i.e. ill-conditioned
 $X^T X$ for a couple reasons

Ex. $X^T X$ not invertible b/c one var is a
LC of the others

Ex. one var is approx. a LC of others

Ex. $P > N$ then $X^T X$ is not invertible
and $K(X^T X) = \infty$

E.g. X measures gene expression in $N=30$
patients across $P=20,000$ genes.

How do we deal with this?

Today: ① Variable selection.

Next time: ② Shrinkage (Ridge/LASSO)

③ Dimensionality Reduction (PCA)

Variable Selection

Goal: pick a subset of important variables and just use this subset.

Q: How do I define "important"?

Two approaches:

① use an individual metric for each var.
and use subset w/ best metric

e.g. p-values for each var.

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

② Calc. a metric for groups of vars
and use group w/ best metric.

Careful: Looking at training metrics can be misleading.

e.g. $RSS_{\text{train}} \downarrow$ as $P \uparrow$

- Solns:
- ① use some test metric (e.g. x-val.)
 - ② use a penalized in-sample metric

→ Ex. Adjusted R^2

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

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

Ex. Mallows's C_p

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

penalty increases w/ p

Ex. AIC

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

Ex. BIC

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

