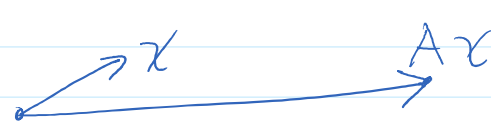


Condition Number (of a matrix)Rel. stability of solving $Az = b$ A is $N \times N$ want to find z .

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


max. amt
that A stretches
a vectorwhy $\|x\|$ a number

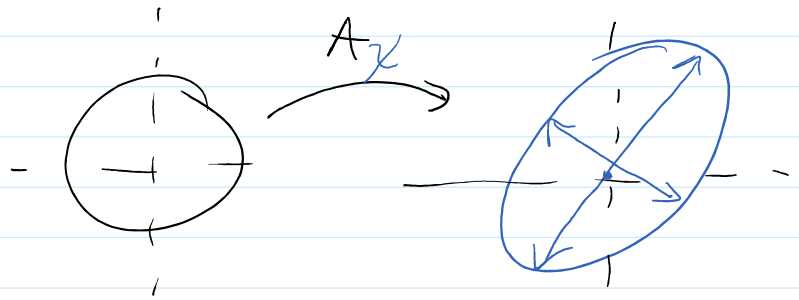
$$\frac{\|Ax\|}{\|x\|} = \|A \underbrace{x/\|x\|}_{y \rightarrow \|y\|=1}\|$$

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

min amt A shrinks any vectornotice: if A is singular then for some x we have $Ax = 0$ so $m = 0 \Leftrightarrow A$ singular.Defn: Condition Number

$$\kappa(A) = M/m = \frac{\text{max stretchy}}{\text{min stretchy}}$$

$$K(A) = M/m = \frac{\text{max stretching}}{\text{min stretching}}$$



Imagine $Az = b$ and $A\delta_z = \delta_b$

then

$$A(z + \delta_z) = b + \delta_b$$

By defn
 $\forall x \quad \|Ax\| \leq M\|x\|$
 and $\|Ax\| \geq m\|x\|$

So $\|b\| = \|Az\| \leq M\|z\|$ (1)

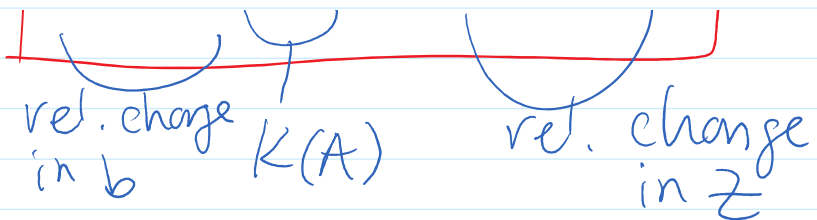
and $\|\delta_b\| = \|A\delta_z\| \geq m\|\delta_z\|$ (2)

re-arrange

by 2 $\frac{\|\delta_b\|}{m} \geq \|\delta_z\|$

multiply:

$$\frac{\|\delta_b\|}{\|b\|} \frac{M}{m} \geq \frac{\|\delta_z\|}{\|z\|}$$



interp: rel change in $b \rightsquigarrow K(A)$
to change in z .

If $K(A)$ is small then system is stable
b/c rel. large change in b
propagate to small changes in soln.

If $K(A)$ is large then small changes to
 b lead to LARGE changes
in soln.

If A is singular (not invertible)
then $K(A) = \frac{M}{m} \leftarrow m=0 = \infty$

in which case $Az = b$
(∞ change in z — no change in b).

Relate back to SVD

Note: If Q is orthogonal

then $\|Qx\| = \|x\|$
(interp: orthog. transf.)

\rightarrow pf.

$$\begin{aligned}\|Qx\|^2 &= (Qx)^T (Qx) \\ &= x^T Q^T Q x \\ &= x^T x\end{aligned}$$

rotate space - don't stretch

$$= \|x'\|^2$$

Let $A = UDV^T$

then

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

$$= \max_{\|x\|=1} \|D \underbrace{V^T x}_y\| \text{ where } \|y\|=1$$

$$= \max_{\|y\|=1} \|Dy\| \rightarrow \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_N)$$

$$= \max_{\|y\|=1} \sqrt{\underbrace{(\sigma_1 y_1)^2}_1 + \underbrace{(\sigma_2 y_2)^2}_0 + \dots + \underbrace{(\sigma_N y_N)^2}_0}$$

set $y_1 = 1$
others to 0 $\sigma_1 \geq \sigma_2 \geq \sigma_3 \geq \dots \geq \sigma_N$

$$= \sigma_1$$

Similarly $m = \sigma_N$

hence $\boxed{K(A) = \sigma_{\max} / \sigma_{\min}}$

Why do we care?

Recall for LS regression, we said

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

this follows b/c if $\frac{\partial L}{\partial \beta} = 0$ ← squared loss
 this leads to system of eqns

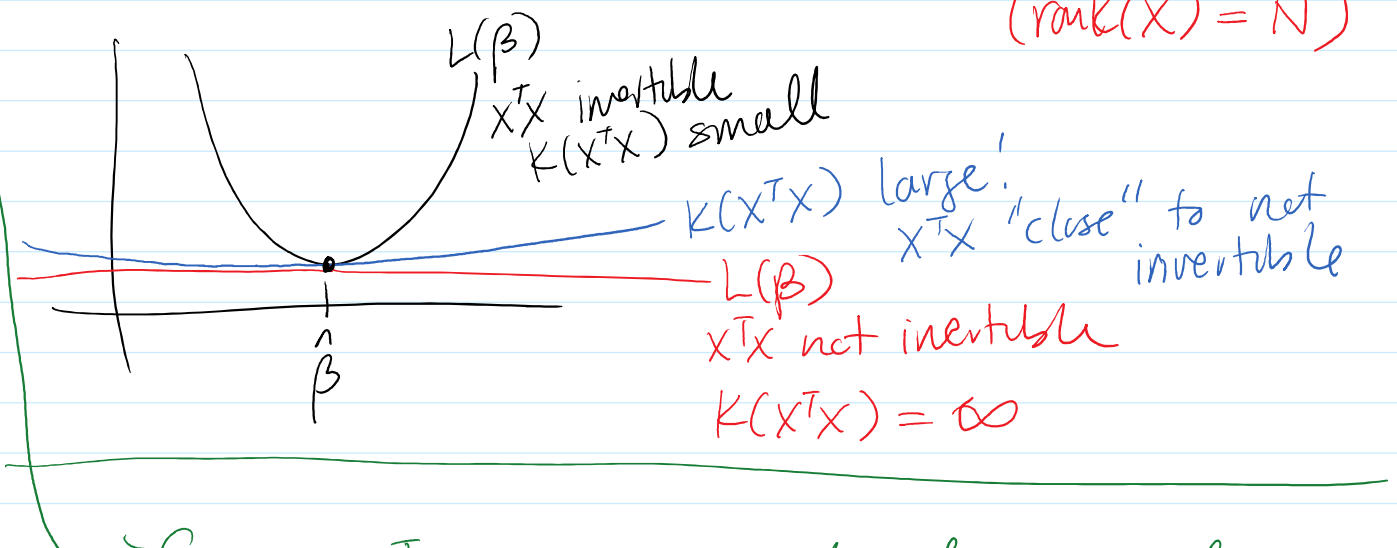
$$\underbrace{(X^T X)}_A \underbrace{\beta}_z = \underbrace{X^T Y}_b$$

← Normal equations

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

Q: In reality when do we get an "ill-conditioned" regression problem i.e. $K(X^T X)$ large?

We said previously that $\hat{\beta}$ unique $\Leftrightarrow X^T X$ invertible
(rank(X) = N)

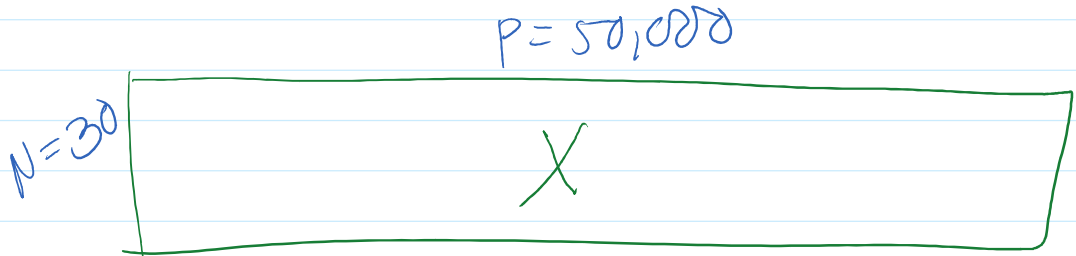


→ Ex. $X^T X$ not invertible if any cols of X are a LC of any other

Ex. more likely $X^T X$ close to not invertible
one variable \approx LC of others

Ex. If $P > N$ then $X^T X$ is not invertible.

e.g. X is gene expression for $P=50,000$
genes of $N=30$ people.



How do we deal w/ this?

- (*) (1) Variable selection
(remove variables to better condition our problem)
- (2) Shrinkage (Ridge / Lasso)
- (3) Dimensionality Reduction (PCA regression)

③ Dimensionality Reduction (PCA regression)

→ unsupervised analysis.

Variable Selection

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

Q: How do we define "important" variables?
→ How do I define a good set of vars?

Two approaches:

① individual metrics

e.g. p-value for each variable, choose vars w/ smallest p-values

problem: performance of one var might be affected by inclusion/exclusion of others

② metrics for sets of variables

idea: build a model w/ different collections of vars and use "best" performing group.

Careful: don't look at trainy metrics
b/c trainy metrics \uparrow as $p \uparrow$
(more flex as $p \uparrow$)

Solve: ① train/val/test split.

② penalize train metrics by p .
(classic)

Ex. Adjusted R^2 ($\uparrow R^2 = \text{better}$)

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

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

RSS-based metrics ($\downarrow \text{RSS} = \text{better}$)

① Mallow's C_p

$$C_p = \frac{1}{N} (\text{RSS} + \underbrace{2p \hat{\sigma}^2}_{\substack{\text{increase w/ } p \\ \text{penalty}}})$$

$\hat{\sigma}^2$ est. of error var

② AIC: Akaike's Information Criterion

$$AIC = \frac{1}{N\hat{\sigma}^2} (RSS + \underbrace{2p\hat{\sigma}^2}_{\text{penalty for large } p})$$

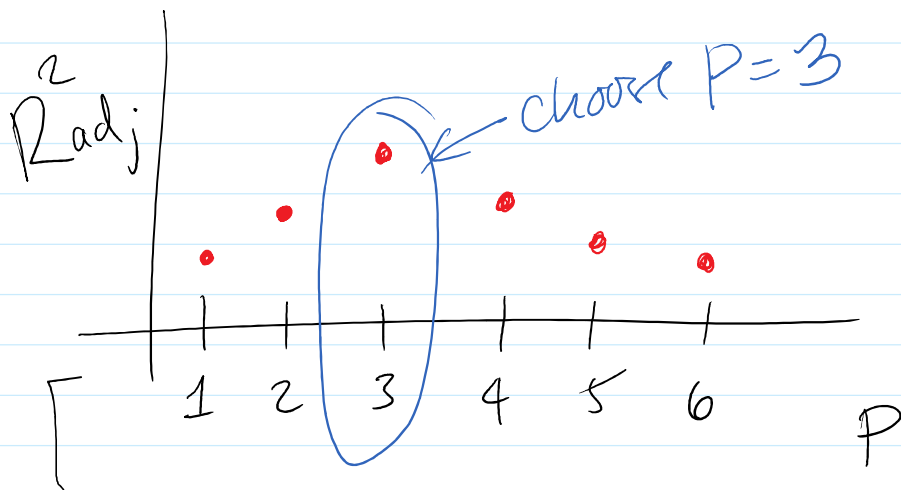
③ BIC: Bayesian Inform. Criteria

$$BIC = \frac{1}{N} (RSS + \underbrace{\log(N)p\hat{\sigma}^2}_{\text{penalty}})$$

Systematic Search: (greedy heuristics)

① Forward ^{stepwise} Selection

- start w/ no covs
- add variables 1 at a time to regr. model
- add var that gives best (inc/dec) in metric



② Backward stepwise Selection

- start w/ all covs.

→ remove one at a time
retrain vari that hurts model metric
the least.

Ideally: test all subsets of vars.
(not computationally possible for large P)
