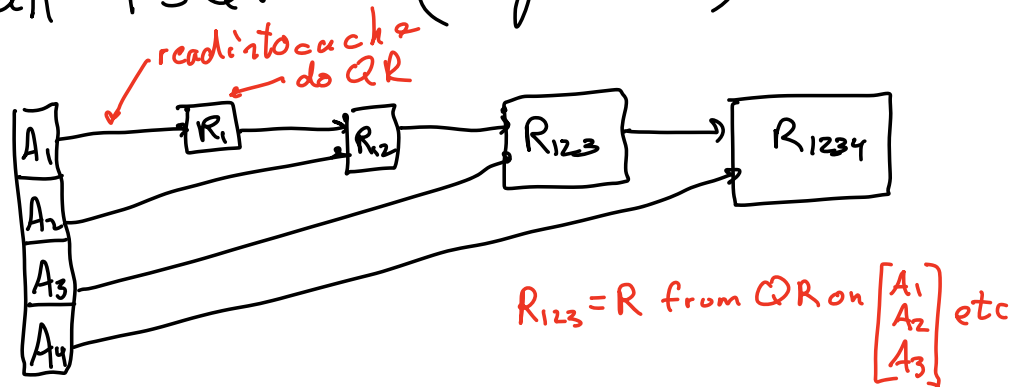
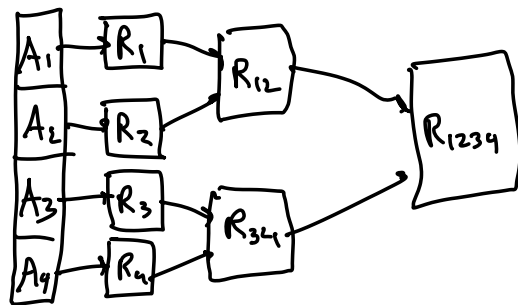


Welcome to Ma221! (Mar 1)

Recall TSQR (sequential)



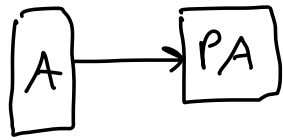
Same idea for parallel TSQR



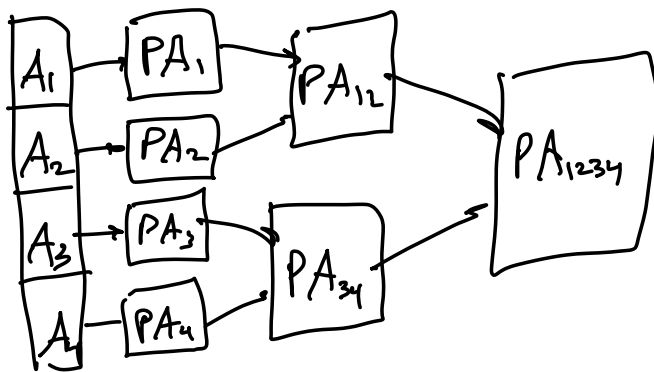
"Map Reduce" where QR is reduction operation

Same idea for partial pivoting on a Tall Skinny matrix (TSLU) one communication tree for many columns

Basic operation



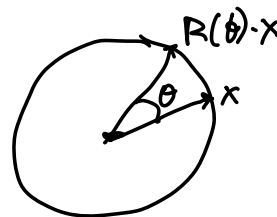
select subset of rows
of A chosen by
partial pivoting
"most linearly independent
rows of A "



Givens rotations: another
simple orthogonal transformation:
used for eigenproblems, SVD

$$R(\theta) = \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix}$$

$R(\theta)x$



$R(i, j, \theta)$ applies rotation

to x_i and x_j

$$\begin{matrix} & & i & & j \\ & & & & \\ i & & \ddots & & \\ & & \cos\theta & -\sin\theta & \\ j & & \sin\theta & \cos\theta & \\ & & & & \ddots \end{matrix}$$

How to pick θ to zero out

$$x_j := \begin{bmatrix} \cos\theta & -\sin\theta \\ \sin\theta & \cos\theta \end{bmatrix} \begin{bmatrix} x_i \\ x_j \end{bmatrix} = \begin{bmatrix} \sqrt{x_i^2 + x_j^2} \\ 0 \end{bmatrix}$$

$$\Rightarrow \cos\theta = \frac{x_i}{\sqrt{x_i^2 + x_j^2}}, \sin\theta = \frac{-x_j}{\sqrt{x_i^2 + x_j^2}}$$

can do QR using Givens rotations,
no advantage over Householder for
dense A , maybe less fillin for
sparse A

Stability of applying Orthogonal
matrices

Summary: Any algorithm that just
multiplies by a sequence of orthogonal
matrices is backward stable

Proof sketch: use basic rule
 $fl(a \text{ op } b) = (a \text{ op } b)(1 + \delta) \quad |\delta| \leq \epsilon$
 to show that multiplying by one
 Householder or Givens **gets** smaller error:

$$fl(Q' \cdot A) = Q' \cdot A + E \quad \|E\| = O(\epsilon) \|A\|$$

Q' "nearly orthogonal" \Rightarrow

$$Q' = Q + F, \quad \|F\| = O(\epsilon)$$

Q exactly orthogonal

$$\begin{aligned} fl(Q' \cdot A) &= Q' \cdot A + E = (Q + F)A + E \\ &= QA + FA + E = QA + G \\ &= \text{exact orthogonal transformation} + G \end{aligned}$$

$$\begin{aligned} \|G\| &= \|FA + E\| \leq \|F\| \cdot \|A\| + \|E\| \\ &\leq O(\epsilon) \|A\| + O(\epsilon) \cdot \|A\| = O(\epsilon) \cdot \|A\| \end{aligned}$$

$$\begin{aligned} fl(Q' \cdot A) &= QA + G = Q(A + Q^T G) \\ &= Q(A + G') \end{aligned}$$

$$\|G'\| = \|G\|$$

\Rightarrow multiplication by Q' is backward stable
 $=$ exact transform of $A + G'$

What if multiply by many Q_i' ?

$$\begin{aligned}
& fl(Q_3'(Q_2'(Q_1'A))) \\
&= fl(Q_3'(Q_2'(Q_1(A + E_1)))) \\
&= fl(Q_3'(Q_2(Q_1(A + E_1) + E_2))) \\
&= (Q_3(Q_2(Q_1(A + E_1) + E_2) + E_3)) \\
&= \underbrace{Q_3 Q_2 Q_1}_Q A + \underbrace{Q_3 Q_2 E_1 + Q_3 E_2 + E_3}_E \\
& \qquad \qquad \qquad \|E\| \leq \|E_1\| + \|E_2\| + \|E_3\| = O(\epsilon) \|A\| \\
&= Q(A + Q^T E) \qquad \|Q^T E\| = \|E\| = O(\epsilon) \|A\|
\end{aligned}$$

same idea (pick $A=I$) to show that
product $fl(Q_3' Q_2' Q_1')$ nearly orthogonal

One more fast (but unstable)

QR algorithm: used in practice
when need Q , know A well conditioned

Cholesky QR:

Factor $A^T A = R^T R$ using Cholesky
form $Q = A \cdot R^{-1}$ since $A = QR$

can fail completely eg if
some pivot during Cholesky is
negative \rightarrow fails, can't do twice

Dealing with (nearly) low rank matrices

Motivation: Real data often low rank (nearly redundant)

- (1) take precautions to avoid inaccurate LS
- (2) use it to compress data, go faster, both deterministically and randomized

Use LS as example of compression, but many applications

Solving a LS problem when A rank deficient

Thm: A $m \times n$, $m \geq n$, rank $r < n$

$$A = U \Sigma V^T = \begin{matrix} & \begin{matrix} r & n-r & m-n \end{matrix} \\ \begin{matrix} m \\ m \end{matrix} & \begin{bmatrix} U_1 & U_2 & U_3 \end{bmatrix} \end{matrix} \begin{matrix} & \begin{matrix} r & n-r \end{matrix} \\ \begin{matrix} r \\ n-r \\ m-n \end{matrix} & \begin{bmatrix} \Sigma_1 & 0 \\ 0 & \Sigma_2 \\ 0 & 0 \end{bmatrix} \end{matrix} \begin{matrix} & \begin{matrix} r & n-r \end{matrix} \\ \begin{matrix} n \\ n \end{matrix} & \begin{bmatrix} V_1 & V_2 \end{bmatrix}^T \end{matrix}$$

Σ_1 is full rank
 $\Sigma_2 = 0$ (later: tiny)

The set of vectors minimizing $\|Ax - b\|_2$ is

$$\left\{ x = V_1 \Sigma_1^{-1} U_1^T b + V_2 y_2, \text{ any } y_2 \in \mathbb{R}^{n-r} \right\}$$

Unique x minimizing both $\|Ax - b\|_2$ and $\|x\|_2$ is gotten from $y_2 = 0$
 $x = V_1 \Sigma_1^{-1} U_1^T b$

Def: $A^+ = V_1 \Sigma_1^{-1} U_1^T$ is Moore-Penrose pseudoinverse of A (includes full rank case, $r = n$)

(in practice: Σ_2 will be all singular values less than some user defined tolerance)

So square or not, full rank or not
"best" solution is $x = A^+ b$

Proof: $\|Ax - b\|_2 = \|U \Sigma V^T x - b\|_2$
 $= \|\Sigma V^T x - U^T b\|_2$ since U orthogonal

$$= \|\Sigma y - U^T b\|_2 \text{ where } y = V^T x$$

$\|x\|_2 = \|y\|_2$ so okay to minimize either one

$$= \left\| \begin{bmatrix} \Sigma_1 y_1 - U_1^T b \\ -U_2^T b \\ -U_3^T b \end{bmatrix} \right\|_2 \text{ where } y = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}$$

minimized by $y_1 = \Sigma_1^{-1} U_1^T b$

$$\text{and } \|y\|_2^2 = \|y_1\|_2^2 + \|y_2\|_2^2$$

minimized by $y_2 = 0$

$$x = Vy = [V_1, V_2] \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = V_1 y_1 + V_2 y_2$$

$$= V_1 \Sigma_1^{-1} U_1^T b$$

Solving LS when A (nearly)
rank deficient, using truncated SVD

$$\kappa(A) = \frac{\sigma_{\max}}{\sigma_{\min}} = \infty \text{ if } A \text{ rank deficient}$$

$$\operatorname{argmin}_x \left\| \begin{bmatrix} 1 & 0 \\ 0 & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} - \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \right\|_2 = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$\operatorname{argmin}_x \left\| \begin{bmatrix} 1 & 0 \\ 0 & \epsilon \\ 0 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} - \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \right\|_2 = \begin{bmatrix} 1 \\ 1/\epsilon \end{bmatrix} \text{ e tiny}$$

what does a "solution" mean if
it can change discontinuously?

Often A not known exactly, just
up to some tolerance, $\|A - A'\|_2 \leq \tau$

What to do?

Def: truncated SVD

$$\text{if } A = U \Sigma V^T$$

$$A(\text{tol}) = U \cdot \Sigma(\text{tol}) \cdot V^T$$

$$\Sigma(\text{tol}) = \text{diag}(\sigma_1(\text{tol}), \sigma_2(\text{tol}), \dots, \sigma_n(\text{tol}))$$

$$\sigma_i(\text{tol}) = \begin{cases} \sigma_i & \text{if } \sigma_i \geq \text{tol} \\ 0 & \text{if } \sigma_i < \text{tol} \end{cases}$$

$A(\text{tol})$ = lowest rank matrix within distance tol of A

Using $A(\text{tol})$ for LS reduces

$$k = \frac{\sigma_{\max}}{\sigma_{\min}} \quad \text{to} \quad \frac{\sigma_{\max}}{\text{tol}}$$

tol is a "knob" for user to trade off sensitivity k and how well LS problem can be solved (how small you can make residual)

Replacing A by an "easier" matrix called regularization, several mechanisms

Lemma: $x_1 = \operatorname{argmin}_x \|A(\text{tol})x - b_1\|_2$
 $x_2 = \operatorname{argmin}_x \|A(\text{tol})x - b_2\|_2$

choose $\|x_i\|_2$ of smallest norm

Then $\|x_1 - x_2\| \leq \frac{\|b_1 - b_2\|}{\text{tol}}$

proof: $\|x_1 - x_2\|_2 = \|(A(\text{tol}))^+ (b_1 - b_2)\|_2$
 $= \|\Sigma(\text{tol})^+ U^T (b_1 - b_2)\|_2$
 $= \|\operatorname{diag}(\sigma_1^{-1}, \sigma_2^{-1}, \dots, \sigma_k^{-1}, 0 \dots 0) \cdot U^T (b_1 - b_2)\|_2$
 \uparrow (last nonzero σ_k where $\sigma_k \geq \text{tol}$)
 $\leq \frac{1}{\sigma_k} \|U^T (b_1 - b_2)\|_2$
 $\leq \frac{1}{\text{tol}} \|b_1 - b_2\|_2$

How does $A(\text{tol})$ depend on tol ?
 piecewise constant, changes
 when $\text{tol} = \sigma_i$

Other advantage of $A(\text{tol})$:
 setting more σ_i to 0, compresses A
 $A(\text{tol})$ rank $k \Rightarrow$ need $m \cdot k + k +$
 $n \cdot k$ words
 to store SVD, can be $\ll m \cdot n$

Solving (nearly) low rank LS
using Tikhonov regularization,
or ridge regression

Replace $\operatorname{argmin}_x \|Ax - b\|_2^2$

by $\operatorname{argmin}_x \|Ax - b\|_2^2 + \lambda \|x\|_2^2$
 $\lambda > 0$

λ "penalizes" very large x
 λ user parameter

$$\operatorname{argmin}_x \|Ax - b\|_2^2 + \lambda \|x\|_2^2$$
$$= \operatorname{argmin}_x \left\| \begin{bmatrix} A \\ \sqrt{\lambda} I \end{bmatrix} x - \begin{bmatrix} b \\ 0 \end{bmatrix} \right\|_2^2$$

\nearrow full rank for any $\lambda > 0$

$$NE \Rightarrow x = \left(\begin{bmatrix} A \\ \sqrt{\lambda} I \end{bmatrix}^T \begin{bmatrix} A \\ \sqrt{\lambda} I \end{bmatrix} \right)^{-1} \begin{bmatrix} A \\ \sqrt{\lambda} I \end{bmatrix}^T \begin{bmatrix} b \\ 0 \end{bmatrix}$$

$$(*) = (A^T A + \lambda I)^{-1} A^T b$$

just add λ to diagonal of $A^T b$

How does λ change SVD?

plug $A = U \Sigma V^T$ into (*)

$$x = V (\Sigma (\Sigma^2 + \lambda I)^{-1}) U^T b$$

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

usual SVD if $\lambda = 0$

$$\sigma_i \gg \lambda^{1/2} \Rightarrow \frac{\sigma_i}{\sigma_i^2 + \lambda} \sim \frac{1}{\sigma_i}$$

$$\sigma_i < \lambda^{1/2} \Rightarrow \frac{\sigma_i}{\sigma_i^2 + \lambda} \leq \frac{1}{\lambda^{1/2}}$$

i.e. λ and tol in $A(\text{tol})$
play analogous roles

Solving low rank LS using QR

QR with column pivoting

Suppose we did $A = QR$ exactly
 $\text{rank}(A) = r < n$, what would R
look like?

If **leading** r columns of A were
linearly independent (true for
"almost all" low rank A)

$$R = \begin{matrix} r & n-r \\ R_{11} & R_{12} \\ n-r \\ 0 & 0 \end{matrix} \quad R_{11} \text{ full rank} \\ \text{so } R_{22} = 0$$

If A nearly low rank, hope that $\|R_{22}\| < \text{tol}$, set $R_{22} = 0$

Assuming this works, solve LS:

$$A = \begin{matrix} n & m \\ Q & Q' \end{matrix} \begin{matrix} R \\ 0 \end{matrix}$$

$$\begin{matrix} n & m \\ Q & Q' \end{matrix} = \begin{matrix} r & n-r & m-n \\ Q_1 & Q_2 & Q' \end{matrix}$$

$$\operatorname{argmin}_x \|Ax - b\|_2$$

$$= \operatorname{argmin}_x \left\| \begin{bmatrix} Q_1 & Q_2 & Q' \\ R \\ 0 \end{bmatrix} x - b \right\|_2$$

$$= \operatorname{argmin}_x \left\| \begin{bmatrix} R_{11} & R_{12} \\ 0 & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} - \begin{bmatrix} Q_1^T b \\ Q_2^T b \\ Q'^T b \end{bmatrix} \right\|_2$$

$$= \operatorname{argmin}_x \left\| \begin{bmatrix} R_{11}x_1 + R_{12}x_2 - Q_1^T b \\ -Q_2^T b \\ -Q_1^T b \end{bmatrix} \right\|_2$$

solution $x_1 = R_{11}^{-1} Q_1^T b - R_{11}^{-1} R_{12} x_2$
for any x_2

How to pick x_2 to minimize $\|x\|_2$?