# Synopsis of Numerical Linear Algebra

Eric de Sturler
Department of Mathematics, Virginia Tech
✉ sturler@vt.edu
🕸 http://www.math.vt.edu/people/sturler

Iterative Methods for Linear Systems: Basics to Research
Numerical Analysis and Software I

$$Ax = b \qquad A \in \mathbb{C}^{n \times n} \quad b \in \mathbb{C}^n \quad \left( \text{or } A \in \mathbb{R}^{n \times n}, \, b \in \mathbb{R}^n \right)$$

$A$ nonsingular (regular) if its column vectors are independent.

   alternatives: $A^{-1}$ exists, $\det(A) \neq 0$, $Ax \neq 0$ for all $x \neq 0$
   
   etc

Otherwise $A$ singular $\left( \text{so } \exists x \neq 0 : Ax = 0, \text{ etc} \right)$

If $A$ nonsingular $A$ has unique solution for any $b \in \mathbb{C}^n$
   $x = A^{-1} b$

If $A$ singular $Ax = b$ has either no solution    or
                                          infinitely many sol.s

   if $Ax = b$ and $Az = 0$
   then $A(x + \alpha z) = b$ for any $\alpha$

In practice, matrix may be nearly singular and this may have an effect on the accuracy of an approximate solution! (more later)

For small (dense) linear systems

LU - decomposition or
Cholesky - decomp.   (also factorization)
and variants

$$A = L \cdot U \quad \begin{cases} L \text{ lower triangular (unit diagonal)} \text{ usually w.} \\ U \text{ upper triangular} \end{cases}$$

$$A x = b \rightarrow L U x = b \quad \begin{cases} \text{Solve } L y = b \text{ forward substitution} \\ \text{Solve } U x = y \text{ backward subst.} \end{cases}$$

$$A \text{ HPD} \quad (A = A^* \text{ and } z^* A z > 0 \text{ for } z \neq 0)$$

$$A = L L^* \quad (\text{or } A = R^* R)$$

For sparse matrices (even quite large ones) special methods
have been developed that reduce fill-in (to reduce complexity
and storage) $\rightarrow$ sparse direct methods

For general matrices pivoting (swapping rows and columns)
is often necessary:

$$P A Q = L U \quad \begin{cases} P \text{ row exchanges} \\ Q \text{ column exchanges} \end{cases}$$
$$(P \text{ or } Q = I \text{ possible})$$

We can also solve linear system using QR decomposition:

$$A^{n \times n} = Q^{n \times n} \cdot R^{n \times n} \quad \begin{cases} Q \text{ is unitary / orthogonal} \\ \\ R \text{ is upper triangular} \end{cases}$$

$$Q^* Q = Q Q^* = I \text{ , } Q \text{ is unitary for complex } Q$$
$$Q^T Q = Q Q^T = I \text{ , } Q \text{ is orthogonal for real } Q$$

Other special matrices

# Other special matrices

$AA^* = A^*A \rightarrow A$ normal (orthog. eigenvectors)

$A = A^* \rightarrow A$ Hermitian

$A = A^T$ (real A) $\rightarrow A$ symmetric

(complex A) $\rightarrow A$ complex symmetric

$x^*Ax > 0$ for any $x \neq 0$, A positive definite
(some people define PD only for symm/Herm matrices, others for all matrices)

# Norms

A norm on a vector space $V$ is any function $f : V \to \mathbb{R}$ such that

1. $f(x) \geq 0$ and $f(x) = 0 \Leftrightarrow x = 0$,
2. $f(\alpha x) = |\alpha| f(x)$,
3. $f(x + y) \leq f(x) + f(y)$,

where $x \in V$ and $\alpha \in \mathbb{R}$.

Important vector spaces in this course: $\mathbb{R}^n$, $\mathbb{C}^n$, and $\mathbb{R}^{m \times n}$, $\mathbb{C}^{m \times n}$ (matrices). Note that the set of all m-by-n matrices (real or complex) is a vector space.

Many matrix norms possess the submultiplicative or consistency property:
$f(AB) \leq f(A) f(B)$ for all $A \in \mathbb{C}^{m \times k}$ and $B \in \mathbb{C}^{k \times n}$ (or real matrices).

Note that strictly speaking this is a property of a *family of norms*, because in general 'each' $f$ is defined on a different vector space.

# Norms

We can define a matrix norm using a vector norm (an induced matrix norm):

$$\| A \|_\alpha = \max_{x \neq 0} \frac{\| Ax \|_\alpha}{\| x \|_\alpha} = \max_{\|x\|_\alpha = 1} \| Ax \|_\alpha$$

Induced norms are always consistent (satisfy consistency property).

Two norms $\| . \|_\alpha$ and $\| . \|_\beta$ are equivalent if there exist positive, real constants $a$ and $b$ such that

$$\forall x : a \| x \|_\alpha \leq \| x \|_\beta \leq b \| x \|_\alpha$$

The constants depend on the two norms but not on $x$.

All norms on a finite dimensional vector space are equivalent.

# Norms

Some useful norms on $\mathbb{R}^n$, $\mathbb{C}^n$, $\mathbb{R}^{m \times n}$, $\mathbb{C}^{m \times n}$:

p-norms: $\|x\|_p = \left( \sum_{i=1}^{n} |x_i|^p \right)^{\frac{1}{p}}$, especially $p = 1, 2, \infty$, where $\|x\|_\infty = \max_i |x_i|$.

Induced matrix p-norms are:

$$\|A\|_1 = \max_j \sum_{i=1}^{n} |a_{ij}| \qquad \text{(max absolute column sum)}$$

$$\|A\|_2 = \sigma_{\max}(A) \quad \text{(max singular value – harder to compute than others)}$$

$$\|A\|_\infty = \max_i \sum_{j=1}^{n} |a_{ij}| \qquad \text{(max absolute row sum)}$$

Matrix Frobenius norm:

$$\|A\|_F = \left( \sum_{i,j=1}^{n} |a_{ij}|^2 \right)^{\frac{1}{2}} \qquad \text{(similar to vector 2-norm for a matrix)}$$

All these norms are consistent (satisfy the submultiplicative property)

## Norms and Inner Products

Let $V$ be a vector space ($\mathbb{C}^n$ or $\mathbb{R}^n$) over complex or
$\quad x, y \in V \qquad \alpha \in \mathbb{C}$ (or $\mathbb{R}$) $\qquad\qquad$ real field

A function $f: V \to \mathbb{R}$ is a norm if

$$f: V \to \mathbb{R} \qquad s.t.$$

(1) $f(x) \geqslant 0$ $\underline{and}$ $f(x) = 0 \iff x = 0$

(2) $f(\alpha x) = |\alpha| f(x)$

(3) $f(x+y) \leq f(x) + f(y)$ $\qquad$ (triangle ineq.)

Often written $\|x\|$, $\|x\|_\alpha$, $\|\|x\|\|$

examples
$$\|x\| = \sum_{i=1}^{n} |x_i|^2$$

Important consequence of triangle ineq. :

$$\Big|\, \|x\| - \|y\| \,\Big| \leq \|x-y\|$$

$$\|x\| = \|y + x - y\| \leq \|y\| + \|x-y\| \iff \|x\| - \|y\| \leq \|x-y\|$$

same way $\|y\| - \|x\| \leq \|y - x\| = \|x-y\|$

Prove induced matrix norm is consistent

Prove for any nonsingular $B$ and any norm $\| \cdot \|_\alpha$ that $\|Bx\|_\alpha$ is a norm

Prove for $B$ HPD (SPD in real case) that $(x^* B x)^{1/2}$ is a norm

$\|$          (for ) for

$\|Bx\|$ is norm    $x$    any nonsingular $B$

$\|x\|$   norm

$\|Bx\| \geq 0$   and   $\|Bx\| = 0 \implies Bx = 0$

$\qquad\qquad\qquad\qquad$ $B$ nonsing $\to x = 0$

$\|Bx\| = 0 \implies x = 0$    $\left.\begin{array}{c}\\\\\end{array}\right\}$ $\|Bx\| = 0 \iff x = 0$

$x = 0 \implies Bx = 0 \implies \|Bx\| = 0$

$\|B(\alpha x)\| = \|\alpha(Bx)\| = |\alpha|\,\|Bx\|$

$\|B(x+y)\| = \|Bx + By\| \leq \|Bx\| + \|By\|$

---

$\|A\|_\alpha = \max\limits_{x \neq 0} \dfrac{\|Ax\|_\alpha}{\|x\|_\alpha}$    show consistent

$\dfrac{\|Ay\|_\alpha}{\|y\|_\alpha} \leq \|A\|_\alpha \iff \|Ay\|_\alpha \leq \|A\|_\alpha \|y\|_\alpha$

$\|AB\|_\alpha = \max\limits_{x \neq 0} \dfrac{\|ABx\|_\alpha}{\|x\|_\alpha}$

$\|AB\|_\alpha \leq \|A\|_\alpha \|B\|_\alpha$   $\boxed{\text{to show}}$

$= \left(\max\limits_{x \neq 0} \dfrac{\|Ax\|_\alpha}{\|x\|_\alpha}\right) \cdot \left(\max\limits_{y \neq 0} \dfrac{\|By\|_\alpha}{\|y\|_\alpha}\right)$

$\max\limits_{x \neq 0} \dfrac{\|ABx\|_\alpha}{\|x\|_\alpha} = \max\limits_{x \neq 0} \dfrac{\|A(Bx)\|_\alpha}{\|Bx\|_\alpha} \cdot \dfrac{\|Bx\|_\alpha}{\|x\|_\alpha}$

$\max \|Au\| \cdot \qquad\qquad \|Bx\|_\sim$

$$\cdots \qquad \|x\|_\alpha \qquad x \neq 0 \qquad \|Bx\|_\alpha \qquad \|x\|_\alpha$$

$$\leq \frac{\max \|Ay\|_\alpha}{\|y\|_\alpha} \cdot \max \frac{\|Bx\|_\alpha}{\|x\|_\alpha}$$

Show $(x^* Bx)^{1/2}$ norm for $B$ HPD

$x^* Bx > 0$ for $x \neq 0$ (by def of PD)

$x = 0 \implies x^* Bx = 0$

So, $(x^* Bx)^{1/2} \geq 0$ and $(x^* Bx)^{1/2} = 0 \iff x = 0$

$$\left((\alpha x)^* B(\alpha x)\right)^{1/2} = (\alpha^* \alpha)^{1/2} (x^* Bx)^{1/2} = |\alpha|(x^* Bx)^{1/2}$$

$$\left((x+y)^* B(x+y)\right)^{1/2} = \left(x^* Bx + x^* By + \underline{y^* Bx + y^* By}\right)^{1/2} =$$

$$\left(x^* Bx + y^* By + 2\,\mathrm{Re}(x^* By)\right)^{1/2}$$

$$\left((x^* Bx)^{1/2} + (y^* By)^{1/2}\right)^2 = x^* Bx + y^* By + 2(x^* Bx)^{1/2}(y^* By)^{1/2}$$

$$|x^* By| \leq (x^* Bx)^{1/2}(y^* By)^{1/2} \quad \rightarrow \text{Cauchy-Schwartz}$$

We'll see many important uses of norms.
First, use to study accuracy and sensitivity of solutions.

Consider two linear systems with slightly different rhs but same matrix.

$$Ax = b \quad\quad \text{and} \quad A(x + \Delta x) = b + \Delta b$$

How "large" is $\Delta x$ (especially relative to length of $x$)
for small (short) $\Delta b$.
Important because
 a) practical problems involve estimates, measurements, noise/error, ...
 b) every computation involves some error (except symbolic but too expensive)
 c) need to understand how sensitive (reliable) answer is

$$A(x + \Delta x) = Ax + A\Delta x = b + \Delta b \iff \Delta x = A^{-1}\Delta b$$

Consider any (vector) norm and its induced matrix norm

$$\|\Delta x\| = \|A^{-1}\Delta b\| \leq \|A^{-1}\|\,\|\Delta b\|$$

$$\frac{\|\Delta x\|}{\|x\|} \leq \|A^{-1}\|\frac{\|\Delta b\|}{\|x\|} = \|A\|\cdot\|A^{-1}\|\cdot\frac{\|\Delta b\|}{\|A\|\cdot\|x\|}$$

from $b = Ax \implies \|b\| \leq \|A\|\,\|x\|$ we get

$$\frac{\|\Delta x\|}{\|x\|} \leq \|A\|\,\|A^{-1}\|\cdot\frac{\|\Delta b\|}{\|b\|}$$

relative change (error) in $x$ is bounded by

$\underline{\|A\|\cdot\|A^{-1}\|}$ times relative change in $b$
$\downarrow$
$\text{cond}(A)$ or $\varkappa(A)$   condition number of $A$

measures sensitivity of solution (conditioning)
depends on chosen norm.


Similar result holds more generally for changes in
A and/or b

$$(A + \Delta A)(x + \Delta x) = b + \Delta b \quad \text{and} \quad Ax = b$$

$$\cancel{Ax} + A\Delta x + \Delta A x + \underbrace{\Delta A \Delta x} = \cancel{b} + \Delta b \quad (Ax = b)$$

neglect products of small factors

$$A\Delta x = \Delta b - \Delta A x \iff \Delta x = A^{-1}\Delta b - A^{-1}\Delta A x$$

$$\|\Delta x\| \leq \|A^{-1}\| \|\Delta b\| + \|A^{-1}\| \|\Delta A\| \|x\| \qquad \iff$$

$$\frac{\|\Delta x\|}{\|x\|} \leq \|A^{-1}\| \frac{\|\Delta b\|}{\|x\|} + \|A^{-1}\| \|\Delta A\| \qquad =$$

$$\|A\| \|A^{-1}\| \frac{\|\Delta b\|}{\|A\| \|x\|} + \|A\| \|A^{-1}\| \cdot \frac{\|\Delta A\|}{\|A\|} \qquad \Rightarrow$$

$$\frac{\|\Delta x\|}{\|x\|} \leq \text{cond}(A) \underbrace{\left( \frac{\|\Delta b\|}{\|b\|} + \frac{\|\Delta A\|}{\|A\|} \right)}$$

$$\qquad\qquad\qquad \text{total relative change in linear system}$$

Storing/computing system and any numerical computation
will introduce errors of roughly machine precision $\varepsilon_{mach}$

double precision : $\varepsilon_{mach} \sim 10^{-16}$

So, at best $\frac{\|\Delta x\|}{\|x\|} \leq \text{cond}(A)\, \varepsilon_{mach}$

In general, we don't know the error in the solution, because we don't know the exact solution.

Let $\hat{x}$ be the (unknown) exact solution to $Ax = b$
Let $\tilde{x}$ be an approximate solution

residual: $r = b - A\tilde{x}$ (easy to compute)

$r = A\hat{x} - A\tilde{x} = Ae$ (e is error $\hat{x} - \tilde{x}$)

$e = A^{-1}r \implies \|e\| \leq \|A^{-1}\|\|r\|$    so, we can bound the

norm of the error directly in terms of the norm of the residual. In practice, we often only estimate $\|A^{-1}\|$

$$\frac{\|e\|}{\|\tilde{x}\|} \leq \|A\| \cdot \|A^{-1}\| \frac{\|r\|}{\|A\|\|\tilde{x}\|} = cond(A) \frac{\|r\|}{\|A\|\|\tilde{x}\|}$$

If $cond(A)$ not (very) large, a small relative residual implies

a small relative error

Example $cond(A) = 10^6$ and $\dfrac{\|r\|}{\|A\|\|\tilde{x}\|} \leq 10^{-14}$

$\dfrac{\|e\|}{\|\tilde{x}\|} \leq 10^{-8}$

In single precision, we can't expect $\dfrac{\|r\|}{\|A\|\|\tilde{x}\|} \lesssim 10^{-8}$ (in general)

so, we cannot expect relative error to be small.

magnitude of cond nr should be considered relative to precision of machine/arithmetic and relative errors in problem (in A and in b).

precision of machine/arithmetic and relative errors in problem (in A and in b).

In estimating errors we often want "backward error" bounds.

"Forward error" is usual notion of error:

difference between exact and approximate solution $(\hat{x} - \tilde{x})$
(divided by solution for relative error)

Alternatively, consider computed, approximate solution the exact solution of a perturbed problem. Then perturbation is backward error.

It is in general not unique. We want to show (in general) that a small perturbation exist such that computed answer is exact $\longrightarrow$ we solved the problem with small backward error.

In general, this is the best we can expect from a good algorithm. The forward error can still be large, when the problem is sensitive. This is one of the main problems with forward error analysis, a worst case analysis for many algorithms yields a very large bound because the forward error can be large for a sensitive problem (but that doesn't mean the algorithm is bad). For algorithms In some cases, it is easy to see that bound is large only for special instances of the problem. For many alg.s this is, unfortunately, not the case. Backward error analysis allows us to separate the two issues. We can show our algorithm produces an exact solution for slightly perturbed problem.

Well-cond. problem $\longrightarrow$ small change from solution original problem

Ill- cond. problem $\longrightarrow$ possibly large change from sol. original prob.

Solve $Ax = b \longrightarrow$ approx. sol. $\tilde{x}$

Any $(\Delta A, \Delta b)$ s.t. $(A + \Delta A)\tilde{x} = b + \Delta b$ provides B.E. (backw. err.)

We're interested in small instances.

Let $r = b - A\tilde{x}$; then $A\tilde{x} = b - r$. So, residual gives a B.E.

Small relative backw. error if $\frac{\|r\|}{\|b\|}$ is small.

Alternative $\left(A + \frac{r\tilde{x}^*}{\|\tilde{x}\|_2^2}\right)\tilde{x} = A\tilde{x} + r = b - r + r = b$ ($\tilde{x}$ exact sol.)

$$\frac{\left\| r\tilde{x}^* / \|\tilde{x}\|_2^2 \right\|}{\|A\|} = \frac{\|r\|_2 \cdot \|\tilde{x}\|_2}{\|A\| \cdot \|\tilde{x}\|_2} = \frac{\|r\|_2}{\|A\| \cdot \|\tilde{x}\|}$$

Remember $\frac{\|e\|}{\|\tilde{x}\|} \leq \text{cond}(A) \cdot \frac{\|r\|}{\|A\|\,\|\tilde{x}\|}$

So, relative FE (forw. err.) bounded by cond nr. times relative BE.

This holds generally.

Note also $\frac{\|e\|}{\|x\|} \leq \kappa(A) \cdot \frac{\|\Delta b\|}{\|b\|}$.

$\overset{\text{cond}(A)}{}$

# Perturbation Analysis

## (Backward error)

Since we cannot compute exactly we are concerned with effects of errors. (usually relative errors)

Consider comp. $f(x)$    exact: $y = f(x)$
                          comp: $\tilde{y}$

In general we want to know $\|\tilde{y} - y\|$ or $\dfrac{\|\tilde{y}-y\|}{\|y\|}$ (forward error). It turns out that this approach has problems (soon). Instead we assume we computed an exact answer to perturbed problem (input) and assess the effect of that perturbation by perturbation analysis ( backward error$_\xi$ = perturbation, it may not be unique)

computed exactly $f(x+\varepsilon)$    (or $\underline{\underline{f(x+\varepsilon) + \Delta f(x+\varepsilon)}}$)

backward error is $\varepsilon$ ; analyze $|f(x+\varepsilon) - f(x)|$

Example: Compute $u^T v$ · $u, v \in \mathbb{R}^n$
$\hookrightarrow u_i, v_i$ machine numbers

exact answer : $\sum u_i v_i$

floating pt. $fl(u^T v) = \Big( u_1 v_1 (1+\varepsilon_1) + u_2 v_2 (1+\varepsilon_2) \Big)(1+\varepsilon_3) + \cdots$

usually we just write $\Big[ \big( u_1 v_1 (1+\varepsilon) + u_2 v_2 (1+\varepsilon) \big)(1+\varepsilon) + u_3 v_3 (1+\varepsilon) \Big](1+\varepsilon)$

and keep in mind that all "$\varepsilon$" are different

but $|\varepsilon| \leq \varepsilon_M$

$$fl(u^Tv) = u_1v_1(1+\varepsilon)^n + u_2v_2(1+\varepsilon)^n + u_3v_3(1+\varepsilon)^{n-1} + \cdots$$

(of $u_i, v_i$ are not mach. numbers it adds
factor $(1+\varepsilon)^2$ to each term?

$$u_1 \to u_1(1+\varepsilon) \quad u_2 \to u_2(1+\varepsilon) \quad etc.)$$

$$\to (\varepsilon_1 + \varepsilon_2 + \cdots)$$

$$fl(u^Tv) \tilde{\approx} u^Tv + u_1v_1 \varepsilon n\varepsilon + u_2v_2 n\varepsilon + \cdots$$

$$|fl(u^Tv) - u^Tv| \leq n|u|^T|v| \varepsilon$$

* sign of each $\varepsilon$ may differ $\to$ assume all
errors accumulate (worst case)

* assume $\varepsilon$ sufficiently small (relative
to $n$) that we can ignore $\varepsilon^2, \varepsilon^3, \cdots$ terms.

* simplify by ~~replacing~~ taking $n\varepsilon|u_3v_3|$

$\varepsilon. s. o.$ $(n-1)\varepsilon|u_3v_3|$ etc

relative error.
$$\frac{n\varepsilon|u|^T|v|}{|u^Tv|} = \frac{|fl(u^Tv)| \overset{-u^Tv}{}}{|u^Tv|}$$

rel. error can be huge if $|u^Tv|$ ~~is~~ very small
compared with $|u|^T|v|$ even if $n\varepsilon$ still small.

$$\|u\|_\infty = \max_i |u_i| = 1 \qquad \|v\|_\infty = 1$$

but $u^Tv \approx 0$ (orthogonal)

relative

So, there is no bound on (forward) error

We may conclude that "simple" computation
of dot product is therefore unreliable.

Backward error analysis:

$$fl(u^T v) \simeq u^T v + \overline{\text{~~~~~~~~~~~~~~~~}} \cdots$$

$$\simeq n \varepsilon u_1 v_1 + n \varepsilon u_2 v_2 + \cdots$$

$$= u_1 v_1 (1 + n\varepsilon) + u_2 v_2 (1 + n\varepsilon) + u_3 v_3 (1 + n\varepsilon) + \cdots$$

$$\simeq u_1 (1 + \tfrac{1}{2} n\varepsilon) v_1 (1 + \tfrac{1}{2} n\varepsilon) + u_2 (1 + \tfrac{1}{2} n\varepsilon) v_2 (1 + \tfrac{1}{2} n\varepsilon) + \cdots$$

$$\text{(assuming } \tfrac{1}{4} n^2 \varepsilon^2 \text{ negligible)}$$

$$= \tilde{u}_1^T \tilde{v}_1 \quad (\text{exactly})$$

$$\tilde{u}_1 = u_1 + \eta \quad \text{where} \quad |\eta_i| \leq \tfrac{1}{2} n \varepsilon_M$$

$$\text{so} \quad \| \eta \|_\infty \leq \tfrac{1}{2} n \varepsilon_M \quad \text{and} \quad \| u \|_\infty = 1$$

So, floating pt. computation of dot product has
small relative backward error.

What explains the huge relative forward error
if $u^T v \simeq 0$ is the "sensitivity" or "conditioning"
of the ~~exam~~ problem. If $u^T v \simeq 0$ the dot
product is ill-conditioned. If $u^T v$ not very small
the problem is well-conditioned.

Backward error analysis let's us separate
the accumulation of numerical error from

the sensitivity of the problem. So, we combine backw. error anal. with perturbation analysis to obtain bounds on computed answers. If an algorithm produces small (bounded) backward error the accuracy depends on the sensitivity

ill-cond → answer may be poor
(but cannot be helped)

well-cond → answer accurate.

# Inner Products

Many methods to select $z_m$ from the Krylov space are related to projections.

We call $f : S \times S \to \mathbb{R}$ an inner product over the real vector space $S$, if for all vectors $x, y, z$ and scalars $\alpha$,

1. $f(x, x) \geq 0$ and $f(x, x) = 0 \Leftrightarrow x = 0$

2. $f(\alpha x, z) = \alpha f(x, z)$

3. $f(x + y, z) = f(x, z) + f(y, z)$

4. $f(x, z) = f(z, x)$

For a complex inner product, $f : S \times S \to \mathbb{C}$, over a complex vector space $S$ we have instead of property (4): $f(x, z) = \overline{f(z, x)}$.

Inner products are often written as $\langle x, y \rangle$, $(x, y)$, or $\langle x, y \rangle_\alpha$, etc..

We say $x$ and $y$ are orthogonal (w.r.t $\alpha$-IP), $x \perp_\alpha y$ if $\langle x, y \rangle_\alpha = 0$.

# Inner products and Norms

Each inner product defines, or induces, a norm: $\|x\| = \sqrt{\langle x, x \rangle}$. (proof?)

Many norms are induced by inner products, but not all. Those norms that are have additional nice properties (that we'll discuss soon).
An inner product and its induced norm satisfy: $|\langle x, y \rangle| \le \|x\|\|y\|$   (CS ineq)

A norm induced by an inner product satisfies the parallelogram equality:
$$\|x + y\|^2 + \|x - y\|^2 = 2\left(\|x\|^2 + \|y\|^2\right)$$
In this case we can find the inner product from the norm as well:

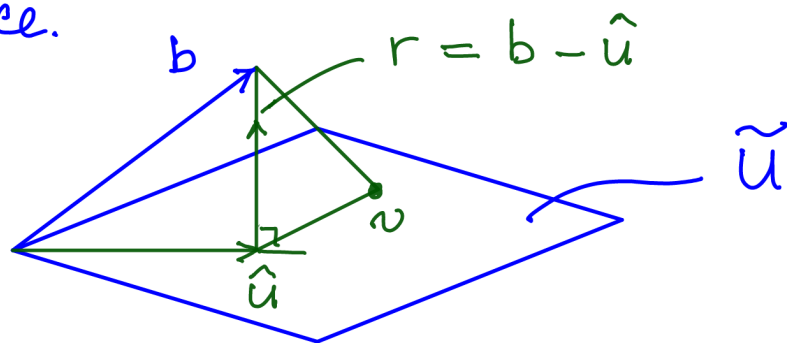Real case:  $\langle x, y \rangle = \dfrac{1}{4}\left(\|x + y\|^2 - \|x - y\|^2\right)$

Complex case:

$$\mathrm{Re}\langle x, y \rangle = \frac{1}{4}\left(\|x + y\|^2 - \|x - y\|^2\right), \quad \mathrm{Im}\langle x, y \rangle = \frac{1}{4}\left(\|x + iy\|^2 - \|x - iy\|^2\right)$$

# Projections

Theo: Best approximation in inner product space.



$$\|b - v\|^2 =$$
$$\|b - \hat{u}\|^2 + \|\hat{u} - v\|^2$$
$$\geq \|b - \hat{u}\|^2$$

Let $\langle \cdot, \cdot \rangle$ be inner prod on $\tilde{V}$ and $\|.\|$ be its associated norm, $\|x\| = \langle x, x \rangle^{1/2}$.
Let $\tilde{U}$ be a subspace of $\tilde{V}$. Then for any $b \in \tilde{V}$ and $\hat{u} \in U$:

$$\|b - \hat{u}\| \leq \|b - v\| \text{ for all } v \in U$$
$$\underline{iff} \quad \langle b - \hat{u}, u \rangle = 0 \quad \text{for all } u \in U.$$

Proof part (i)

Assume $\langle b - \hat{u}, u \rangle = 0$ for all $u \in \tilde{U}$, and $v \in U$.

Then $\|b - v\|^2 = \|b - \hat{u} + \hat{u} - v\|^2 =$

$\langle (b - \hat{u}) + (\hat{u} - v), (b - \hat{u}) + (\hat{u} - v) \rangle =$

$\langle b - \hat{u}, b - \hat{u} \rangle + \langle \hat{u} - v, \hat{u} - v \rangle + \langle b - \hat{u}, \hat{u} - v \rangle + \langle \hat{u} - v, b - \hat{u} \rangle$

Since $\hat{u} - v \in U$, $\langle b - \hat{u}, \hat{u} - v \rangle = \langle \hat{u} - v, b - \hat{u} \rangle = 0.$

So, $\|b - v\|^2 = \|b - \hat{u}\|^2 + \|\hat{u} - v\|^2 \geq \|b - \hat{u}\|^2.$

$$qed$$

It turns out that the orthogonality condition
$$\langle b-\hat{u}, u\rangle = 0 \text{ for all } u \in \tilde{U}$$
provides a recipe for finding $\hat{u}$ that reduces
an optimization problem (hard/expensive)
to solving a linear system of equations
(relatively easy/cheap).

We will show this by construction, which will also
show that $\hat{u}$ is unique for a given $b$.

part (ii)
Assume that $\|b-\hat{u}\| \leq \|b-v\|$ for all $v \in \tilde{U}$.
Now, let $z \in \tilde{U}$, $z \neq 0$, and $\langle b-\hat{u}, z\rangle = t$, and
consider the points $\hat{u}+\alpha t z$ for $\alpha \in \mathbb{R}$.
   (note $\hat{u}+\alpha t z \in \tilde{U}$ for any $\alpha$)
Then $\|b-(\hat{u}+\alpha t z)\|^2 = \langle (b-\hat{u})-\alpha t z, (b-\hat{u})-\alpha t z\rangle$
$= \langle b-\hat{u}, b-\hat{u}\rangle - \alpha t \langle z, b-\hat{u}\rangle - \alpha \bar{t}\langle b-\hat{u}, z\rangle + \alpha^2 |t|^2 \langle z, z\rangle$
$= \|b-\hat{u}\|^2 - 2\alpha |t|^2 + \alpha^2 |t|^2 \|z\|^2$
If $|t|^2 \neq 0$, this is a parabola in $\alpha$ with a
 min $\underline{\text{less than}}$ $\|b-\hat{u}\|^2$.    (just minimize expression)
Hence $\|b-\hat{u}\|^2 \leq \|b-v\|^2$ for all $v \in \tilde{U}$ implies
that $|t|=0 \iff \langle b-\hat{u}, z\rangle = 0$. Since this
must hold for any $z \in \tilde{U}$, $\langle b-\hat{u}, u\rangle = 0$ for all $u \in \tilde{U}$.

Next, given b we construct $\hat{u}$.

Let $u_1, \ldots, u_k$ form a basis for $\tilde{U}$. Then
$\langle b - \hat{u}, u_i \rangle = 0$ for $i = 1..k$
Also $\hat{u} \in U \implies \hat{u} = \sum_{j=1}^{k} u_j \alpha_j$ (where $\alpha_j$ unknown)
Subst. gives:

$$\langle b - \sum_j u_j \alpha_j, u_i \rangle = 0 \qquad i = 1..m \quad \Longleftrightarrow$$
$$\langle b, u_i \rangle - \sum_j \alpha_j \langle u_j, u_i \rangle = 0 \quad \text{for } i = 1 \ldots m$$

Now let $g_{ij} = \langle u_j, u_i \rangle$, $G = (g_{ij})$, the matrix with coeff.s $g_{ij}$ and $f_i = \langle b, u_i \rangle$

Lemma: If $\{u_1, \ldots, u_k\}$ are lin. indep, then the Gram matrix $G$ is invertible. (proof later)

Now we can write $\langle b, u_i \rangle - \sum_j \alpha_j \langle u_j, u_i \rangle = 0$

as $\sum_{j=1}^{m} g_{ij} \alpha_j = b_i$ for $i = 1..m$ \quad or \quad in

matrix form $G\alpha = f \implies \alpha = G^{-1} f$ \quad and

$\hat{u} = \sum_{j=1}^{k} u_j \alpha_j$.

( note that $f$ depends linearly on $b$ and
$\alpha$ depends linearly on $f$. Hence, $\hat{u}$ depends
linearly on $b$)

Best approx. in $\mathbb{C}^m$ with std in.prod / norm

Let $\tilde{u}$ be subspace of $\mathbb{C}^m$ and $\{u_1, ..., u_k\}$ be basis for $\tilde{u}$.

Then $u_i \in \mathbb{C}^m \longrightarrow$ matrix $U = (u_1 \; \cdots \; u_k)$

$$R(U) = \tilde{u} \text{ and } \hat{u} = U\alpha \text{ for some } \alpha \in \mathbb{C}^k$$

Using std. inner product $g_{ij} = \langle u_j, u_i \rangle = u_i^* u_j$ and $G = U^* U$

The orthog. cond. $b - U\alpha \perp \tilde{u} \Rightarrow$

$$b - U\alpha \perp u_i \text{ for } i = 1..k \longrightarrow U^*(b - U\alpha) = 0$$
$$\Longleftrightarrow U^* b - U^* U\alpha = 0 \longrightarrow G\alpha = f \text{ with}$$
$$G = U^* U \text{ and } f = U^* b$$

Again $\alpha = G^{-1} f = (U^* U)^{-1} U^* b$ and
$$\hat{u} = U\alpha = U(U^* U)^{-1} U^* b$$

Proof of Lemma for Gram matrix for this case:
$U^* U$ is invertible iff _____ lin. indep.

Assume $\{u_1, ...u_k\}$ is lin. indep.

Let $U^* U z = 0$. Then $z^* U^* U z = 0 \Longleftrightarrow \|Uz\|_2^2 = 0$, which implies $Uz = 0$. Since the columns of $U$ are lin. indep. $z = 0$ must hold. Hence $U^* U$ is invertible.

Conversely, assume $U^*U$ is invertible.

Let $\displaystyle\sum_{j=1}^{k} u_j z_j = Uz = 0$. Then $U^*Uz = U^*0 = 0$

and $z = (U^*U)^{-1}(U^*U)z = 0$.

The same arguments apply for a general inner product.


We saw that the best approx. problem to $b$ in the space $\tilde{U} \subseteq \mathbb{C}^m$ was solved by $\hat{u} = U(U^*U)^{-1}U^*b$

Let $P = U(U^*U)^{-1}U^*$. Verify that $P^2 = P$. Such a matrix or, more generally, linear transformation is called a projection (or projector).

Def: let $P : \tilde{V} \to \tilde{V}$ be a linear transf. and $P^2 = P$. Then $P$ is a projection.

The particular choice for $P$ above satisfies another important property.

Def: Let $\tilde{V}$ be a vector space w. inner product $\langle \cdot, \cdot \rangle_{\tilde{V}}$. A linear transf. $T: \tilde{V} \to \tilde{V}$ (or linear oper.) is called selfadjoint if

$$\forall x, y \in \tilde{V} : \langle Tx, y \rangle_{\tilde{V}} = \langle x, Ty \rangle_{\tilde{V}}$$

In $\mathbb{C}^m$ with the std inner product $(\langle x, y \rangle = y^* x)$ a matrix is self-adjoint if it is Hermitian:

$$A \in \mathbb{C}^{m \times m} : \quad \langle Ax, y \rangle = \langle x, Ay \rangle \iff$$
$$y^* A x = y^* A^* x$$

Taking $x, y = e_i, e_j$ (Cartesian basis vectors) for $i, j = 1 .. m$ shows this component-wise.

$P = U(U^* U)^{-1} U^*$ is self-adj. wrt the std in.prod.

Def: An orthog. proj. is a lin. transf. $P: \tilde{V} \to \tilde{V}$ such that $P = P^2$ and $P$ is self-adjoint wrt $\langle \cdot, \cdot \rangle_{\tilde{V}}$.

An orthog. proj. $P$ has the property that $\|b - Pb\|_V \leq \|b - u\|_V$ for all $u \in R(P)$ (already proved above)

Def: Let $\tilde{U}, \tilde{V} \subseteq \tilde{W}$ ($\tilde{W}$ vector space)
$$\tilde{U} + \tilde{V} = \{u + v \mid u \in \tilde{U}, v \in \tilde{V}\}$$

<u>Def</u>: Let $\tilde{U}, \tilde{V}$ be subspaces of $\tilde{W}$ (vector space)
We say that $\tilde{U}$ and $\tilde{V}$ form a <u>direct sum</u>
<u>decomposition</u> of $\tilde{W}$: $\tilde{U} \oplus \tilde{V} = \tilde{W}$ if

1) $\tilde{U} + \tilde{V} = \tilde{W}$ (or $\forall w \in \tilde{W}: w = u + v$, where
$\qquad\qquad\qquad\qquad u \in \tilde{U}, v \in \tilde{V}$)

2) $\tilde{U} \cap \tilde{V} = \{\underline{0}\}$

(note two subspaces $\tilde{U}, \tilde{V}$ form a direct sum
of $\tilde{U} \oplus \tilde{V}$ if $\tilde{U} \cap \tilde{V} = \{\underline{0}\}$. )

## Properties of Projections

Let $P: \tilde{V} \to \tilde{V}$ be projection. Then $(I-P)$ is projection (called complementary projection)

$(I-P)^2 = I - P - P + P^2 = I - P$

$R(P) = N(I-P), \quad N(P) = R(I-P) \quad$ (verify)

__Theo:__ $R(P) \oplus N(P) = \tilde{V}$

Proof: let $v \in \tilde{V}$. Then $v = v - Pv + Pv =$

$\quad (I-P)v + Pv$, where $(I-P)v \in N(P)$
$\qquad\qquad\qquad\qquad\qquad\quad Pv \qquad \in R(P)$

So, $R(P) + N(P) = \tilde{V}$

Let $x \in R(P)$ and $x \in N(P)$. Then $x = Py$, for some $y$, and $Px = 0$. So, $0 = Px = P^2 y = Py = x$.
Hence, $R(P) \cap N(P) = \{0\}$.

Note that $x \in R(P) \Rightarrow Px = x$, so $x$ eig. vec. with eig. val 1. Also, $x \in N(P) \Rightarrow Px = 0$, so $x$ is eig. vec with eig. val. 1.

Since every vector $v \in \tilde{V}$ can be written (uniquely) as $v = x+y$, where $x \in R(P)$ and $y \in N(P)$, the union of a basis for $R(P)$ and a basis for $N(P)$ is a basis for $V$ (this holds generally for direct sums). Hence $P$ is diagonalizable and has eigenvalues $\Lambda(P) = \{0, 1\}$.

So, if $P: \mathbb{C}^m \to \mathbb{C}^m$ and $\{r_1, ..., r_k\}$ basis for $R(P)$, $\{n_1, ..., n_\ell\}$ basis for $N(P)$, then

$$P(r_1 ... r_k \ n_1 ... n_\ell) = (r_1 ... r_k \ n_1 ... n_\ell) \begin{pmatrix} I_k & \\ & 0_\ell \end{pmatrix}.$$

For $P = U(U^*U)^{-1}U^*$, $R(P) = R(U) = \tilde{U}$.

$N(P)$?

Since $U$ has lin. indep. col.s (basis for $\tilde{U}$)

$Py = 0$ iff $U^*y = 0 \Rightarrow u_1^*y = 0 ... u_k^*y = 0$

$\quad$ that is $\langle y, u_j \rangle = 0$ for $j = 1..k$ $\Rightarrow$

$\qquad \langle y, u \rangle = 0$ for all $u \in \tilde{U}$

$\quad$ The set of all such $y$ is called the orthog.
complement of $\tilde{U}$:

$$\tilde{U}^\perp = \{ y \in \mathbb{C}^m \mid \langle y, u \rangle = 0 \text{ for all } u \in \tilde{U} \}$$

So, $N(P) = \tilde{U}^\perp$.


So, a projection always defines a direct sum
decomposition:

$$R(P) \oplus N(P) = V$$

but not necessarily $N(P) = R(P)^\perp$. The latter
holds only for orthogonal proj.s

It turns out that, conversely, every direct sum defines a projection.

Let $\tilde{U} \oplus \tilde{V} = \tilde{W}$.

Theo: The sum $w = u + v$, with $u \in \tilde{U}$, $v \in \tilde{V}$ (for any $w \in \tilde{W}$) is unique.

Pr. Let $w \in \tilde{W}$ and $w = u_1 + v_1$, $u_1 \in \tilde{U}$, $v_1 \in \tilde{V}$

$\quad\quad\quad$ and $\quad w = u_2 + v_2 \quad u_2 \in \tilde{U}$, $v_2 \in \tilde{V}$

Then $u_1 + v_1 = u_2 + v_2 \iff u_1 - u_2 = v_2 - v_1$

Since $\tilde{U}$, $\tilde{V}$ are vector spaces $u_1 - u_2 \in \tilde{U}$ and $v_2 - v_1 \in \tilde{V}$. But then $\tilde{U} \cap \tilde{V} = \{\underline{0}\}$ (from def. of direct sum) implies

$u_1 - u_2 = 0 \implies u_1 = u_2$ $\left.\right\}$ Hence the $u, v$ in
$v_2 - v_1 = 0 \implies v_1 = v_2$ $\left.\right\}$ $w = u + v$ are unique.

So, following discussion of projections above define $P$ such that $R(P) = \tilde{U}$ and $N(P) = \tilde{V}$ and $P^2 = P$.

Let $w \in \tilde{W}$ and $w = u + v$ $(u \in \tilde{U}, v \in \tilde{V})$

$N(P) = \tilde{V} \implies Pw = Pu + \underset{=0}{\cancel{Pv}} = Pu$

$P^2 = P \quad\quad \implies Pu = u$

$\quad\quad\quad (u \in R(P) \implies u = Px, \text{ for some } x$
$\quad\quad\quad\quad\quad \implies Pu = P^2 x = Px = u )$

So $P w = u$   (verify that $P$ is linear)
  $\hookrightarrow$ the $\tilde{u}$ component of $w$.

Note that for $v \in \tilde{V}$ : $\underset{\in \tilde{u}}{v = \underbrace{0}} + \underset{\in \tilde{V}}{\underbrace{v}}$   and

$P v = \underline{0} \implies \tilde{V} \subseteq N(P)$

Conversely, if $x \in N(P)$, then $x = \underline{0} + x_V$ , $x_V \in \tilde{V}$
and $x \in \tilde{V}$. So, $N(P) \subseteq \tilde{V}$. Hence, $N(P) = \tilde{V}$.

In practice to find the $\tilde{u}$ comp. and $\tilde{V}$ comp.
(so, we can pick the $\tilde{u}$ comp.) would be
expensive − essentially a large linear solve.
So, we define projections using an auxiliary space
(or basis for such a space) and the complemen-
tary space $\tilde{V}$ is defined implicitly.

Example: Find $\hat{u} \in \tilde{u}$ s.t. $b - \hat{u} \perp \tilde{Y}$
If $\dim \tilde{u} = k$ we need in principle $\dim \tilde{y} = k$:
$b - \sum u_j \alpha_j \perp y_1 \dots y_k \implies$ (std in. prod.)
  $\quad y^*(b - U\alpha) = 0 \iff y^* b = y^* U \alpha$
$\alpha$ uniquely defined for any $b$ if $y^* U$ invertible
Let $P b = U \alpha \implies R(P) = \tilde{u}$.
What is $\tilde{V} = N(P)$ ?

$Pb = \underline{0} \implies \alpha = 0 \implies \alpha = (y^*u)^{-1} y^* b = 0$

$\implies b \in R(y)^\perp$. So, $\tilde{V} = R(y)^\perp$.

$Pb = u(y^*u)^{-1} y^* b \quad \rightarrow \quad P = u(y^*u)^{-1} y^*$

This defines a projection if $y^*u$ invertible which is equivalent to $R(y)^\perp \cap R(u) = \{\underline{0}\}$.

As we compute $\tilde{u}$ (or its basis) and $\tilde{y}$ (or its basis) simultaneously, at
the projection may not be defined and the method breaks down.

If $\tilde{y} = \tilde{u}$, so $\tilde{V} = \tilde{y}^\perp$, $P$ is an orthogonal projection (wrt the std inner product).

# Approximations from subspaces

If we pick an approx. from a subspace and set the residual orthogonal to that subspace, we often call the approx. a (Ritz) Galerkin approx. (and the cond. a Galerkin cond.)
For example, solving approx. $Ax = b$ using $\tilde{U}$ subspace of $\mathbb{C}^m$ to find approx. :

Find $\hat{x} \in \tilde{U}$ s.t. $b - A\hat{x} \perp \tilde{U}$ (in some inner product).
If $R(U) = \tilde{U}$ : $b - AU\xi \perp U$
std in.prod. gives $U^*(b - AU\xi) = 0$.

If we use an additional space $\tilde{V}$ to define the projection, we call the approx. a Petrov-Galerkin approx. (Petrov-Galerkin condition).
Find $\hat{x} \in \tilde{U}$ s.t. $b - A\hat{x} \perp \tilde{V}$ $(\dim \tilde{V} = \dim \tilde{U})$
If $R(V) = \tilde{V}, R(U) = \tilde{U}$ and std. in.prod. used

$$V^*(b - AU\xi) = 0$$

$$A \in \mathbb{C}^{n \times m} \qquad n \geq m$$
$$b \in \mathbb{C}^{n}$$

$$\min_{x \in \mathbb{C}^m} \| b - Ax \|_2 \qquad (\text{also written } Ax \approx b)$$

If $b \in \text{range}(A)$ then $\exists x : Ax = b$ (by def.)

What if $b \notin \text{range}(A)$?



$r = b - Ax \perp \text{range}(A)$

$\text{range}(A)$

$Ax$

$$A^* r = A^*(b - Ax) = 0 \iff A^* b - A^* A x = 0$$

Slight simplification: col.s A independent.

Then $A^* A$ nonsingular (why?): $x = (A^*A)^{-1} A^* b$

$$Ax = A(A^*A)^{-1}A^* b \qquad b - Ax = (I - A(A^*A)^{-1}A^*)b$$

Let $P = A(A^*A)^{-1}A^* b$ $\begin{cases} P^2 = A(A^*A)^{-1} A^*A (A^*A)^{-1} A^* = P \\ \qquad\qquad\qquad \hookrightarrow P \text{ projection} \\ P = P^* \text{ Herm/Symm (if real)} \end{cases}$

$P^2 = P$ <u>and</u> $P = P^* \to P$ orthogonal projection
$\qquad\qquad$ (since $\text{range}(I-P) \perp \text{range}(P)$)

$$P^*(I-P) = P^* - P^* P = P - P^2 = P - P = 0$$

$b = Ib = (P + (I-P))b = Pb + (I-P)b = b_1 + b_2$ $\begin{cases} b_1 \in \text{range}(A) = \text{range}(P) \\ b_2 \perp \text{range}(A) \end{cases}$

Solve $Ax = b_1 \qquad r = (I-P)b = b_2$

Theo: Let $x = \arg\min\limits_{\tilde{x} \in \mathbb{C}^m} \| b - A\tilde{x} \|_2$ .

Then $b - Ax \perp \text{range}(A)$

Proof: Consider $f(x) = \langle b - Ax, b - Ax \rangle$ (quadr. function)

minimum for $\hat{x}$ ( $f(\hat{x})$ minimal )

Then for any $p \in \mathbb{C}^m \neq 0$ we must have

$$\frac{d}{dp} f(\hat{x}) = 0 \quad \Rightarrow \quad \lim_{\substack{t \to 0 \\ (t \in \mathbb{R})}} \frac{f(\hat{x} + tp) - f(\hat{x})}{t} = 0 \quad (\text{any } p \in \mathbb{C}^m)$$

$$\frac{\langle b - A(\hat{x} + tp), b - A(\hat{x} + tp) \rangle - \langle b - A\hat{x}, b - A\hat{x} \rangle}{t} =$$

$$\frac{1}{t} \left( \langle b - A\hat{x}, b - A\hat{x} \rangle - \langle b - A\hat{x}, Atp \rangle - \langle Atp, b - A\hat{x} \rangle + \langle Atp, Atp \rangle - \langle b - A\hat{x}, b - A\hat{x} \rangle \right) =$$

$$\frac{1}{t} \left( -t \langle b - A\hat{x}, Ap \rangle - t \langle Ap, b - A\hat{x} \rangle + t^2 \langle Ap, Ap \rangle \right) =$$

$$- \langle b - A\hat{x}, Ap \rangle - \langle Ap, b - A\hat{x} \rangle + t \langle Ap, Ap \rangle$$

$$\lim_{t \to 0} \; - \langle b - A\hat{x}, Ap \rangle - \langle Ap, b - A\hat{x} \rangle + t \langle Ap, Ap \rangle = 0 \quad (\text{any } p) \quad \Longleftrightarrow$$

$$\langle b - A\hat{x}, Ap \rangle + \langle b - A\hat{x}, Ap \rangle = 0 \quad (\text{any } p) \quad \overset{(\text{why ?})}{\Longleftrightarrow}$$

$$\langle b - A\hat{x}, Ap \rangle = 0 \quad (\text{any } p)$$

$$b - A\hat{x} \perp Ap \text{ for any } p \quad \Longleftrightarrow \quad b - A\hat{x} \perp \text{range}(A)$$

proj. /orthog. proj.

$P^2 = P$ and not $P^* = P$ $\Rightarrow$ oblique /skew proj

$U = \text{range}(P)$ $\quad \rightarrow \quad P : \mathbb{C}^n \rightarrow U$

$\quad P^2 = P \Rightarrow P_U = I_U$ (restr. of $P$ to $U$)

any $u \in U$ : $\exists z$ s.t. $Pz = u$

then ~~$\text{(...)}$~~ $Pu = P(Pz) = P^2 z = Pz = u$

$\Rightarrow$ any proj.

$V = \text{range}(I - P)$ $\qquad P^*(I - P) = P^* - P^* P \neq 0$

$\qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad$ (in gen)

$\quad v \in V \rightarrow v = (I - P) z$

$\qquad \qquad P v = (P - P^2) z = 0$

~~$\text{(...)}$~~ $P$ if $Pv = 0 \Rightarrow (I - P) v = v$

$v \in \text{range}(I - P)$ (also proj. so —)

So $v \in \text{null}(P) \Rightarrow v \in \text{range}(I - P)$, and

$v \in \text{range}(I - P) \Rightarrow v \in \text{null}(P)$, equal

$z \neq 0 \quad z \in \text{range}(P) \Rightarrow z \notin \text{null}(P)$

$\qquad \qquad z \in \text{null}(P) \Rightarrow z \notin \text{range}(P)$

disjoint.

any $z = z_1 + z_2$ s.t.

$z_1 \in \text{range}(P)$ and $z_2 \in \text{null}(P) = \text{range}(I - P)$

$\rightarrow$ direct sum $\begin{cases} z = Iz = (P + (I - P)) z = \\ \quad Pz + (I - P) z = \\ \quad z_1 \quad + \quad z_2 \end{cases}$

given $U, V$ as above $\longrightarrow$ any $z$: $\hat{u}, \hat{v}$

$$z = \hat{U}\zeta_1 + \hat{V}\zeta_2$$

range$(\hat{u}) = u$
range$(\hat{v}) = V$

$$Pz = U\zeta_1 \qquad (I-P)z = V\zeta_2$$

$$P^2 z = PU\zeta_1 = U\zeta_1 \qquad -\ -$$

$$\begin{cases} P = [\hat{u}\ \hat{v}]\begin{bmatrix} I & \\ & 0 \end{bmatrix}[\hat{u}\ \hat{v}]^{-1} \\[4mm] I-P = [\hat{u}\ \hat{v}]\begin{bmatrix} 0 & \\ & I \end{bmatrix}[\hat{u}\hat{v}]^{-1} \end{cases}$$

Solving least squares problems: $A = [a_1 \cdots a_m]$

1) (Modified) Gram-Schmidt

① $q_1 = a_1 / \|a_1\|_2$ and $r_{11} = \|a_1\|_2 \rightarrow a_1 = q_1 r_{11}$

② $r_{12} = q_1^* a_2$ , $r_{22} = \|a_2 - q_1 r_{12}\|_2$

$q_2 = (a_2 - q_1 r_{12})/r_{22} \iff r_{22} q_2 = a_2 - q_1 r_{12} \iff a_2 = [q_1 \ q_2]\begin{bmatrix} r_{12} \\ r_{22} \end{bmatrix}$

$[a_1 \ a_2] = [q_1 \ q_2]\begin{bmatrix} r_{11} & r_{12} \\ 0 & r_{22} \end{bmatrix}$

note $q_i^* q_i = 1 \quad (i = 1,2)$ and

$q_1^* q_2 = (q_1^* a_2 - q_1^* q_1 q_1^* a_2) r_{22}^{-1} = 0 \rightarrow q_2 \perp q_1$

⊛ $r_{ik} = q_i^* a_k \quad i = 1 .. k-1 \quad r_{kk} = \|a_k - \sum_{i=1}^{k-1} q_i r_{ik}\|_2$

$q_k = (a_k - q_1 r_{1k} - q_2 r_{2k} - \cdots - q_{k-1} r_{k-1 k}) r_{kk}^{-1}$

$q_i^* q_k = q_i^* a_k - q_i^* q_i r_{ik} = 0 \quad (q_i^* q_j = \delta_{ij})$

$\hookrightarrow$ prove by induction

similarly prove by induction range$([a_1 \cdots a_k]) = $range$([q_1 \cdots q_k])$

$A^{n \times m} = Q^{n \times m} R^{m \times m}$ where

$Q^* Q = I_m$ and $R = \begin{pmatrix} r_{11} & & \\ & \ddots & \\ 0 & & r_{mm} \end{pmatrix}$

If $\{a_1, \ldots, a_m\}$ indep $\rightarrow \{q_1, \ldots, q_m\}$ indep and

R nonsingular

[If col.s A dependent we can still generate

$AP = Q^{n \times m} R^{m \times m}$ with $Q^* Q = I$ and

$\downarrow \qquad\qquad \mathcal{R} \begin{pmatrix} r_{11} & \\ & \searrow \end{pmatrix}$

$$A\Pi = U \quad UR \quad \text{with} \quad U^*U = \underline{V} \quad \text{and}$$

column exchanges $\quad R = \begin{pmatrix} r_{11} & \searrow & \overline{\phantom{r_{pp}}} \\ & \searrow r_{pp} - r_{pm} \\ O & & O \end{pmatrix}$

with slightly adapted algorithm

for accuracy compute $r_{ik} = q_i^* \left( a_k - \sum_{j=1}^{i-1} q_j r_{jk} \right) \quad i = 1 .. m-1$

also if $\left\| a_k - \sum_{j=1}^{k-1} q_j r_{jk} \right\| \ll \| a_k \|$ then likely

compute $q_k$ inaccurate $(\text{not} \perp q_1 \cdots q_{k-1})$

may need reorthogonalization if $Q^*Q = I$ important

Assume $A$ indep cols

Since $\text{range}(Q) = \text{range}(A)$ and $Q^*Q = \underline{I}_m$ :

$QQ^* = A(A^*A)^{-1}A^* \quad (\text{verify})$
$QQ^*b = b_1 \quad (I - QQ^*)b = b_2$

$Q^*b - Q^*Ax = 0 \iff Q^*b - Rx = 0 \quad (\text{solve})$

$b - Ax = r + QQ^*b - QRx = r + Q(Q^*b - Rx)$

$\| b - Ax \|_2^2 = \| r \|_2^2 + \| Q^*b - Rx \|_2^2$

$\| Qz \|_2 = \left( z^* \cancel{Q^*Q} z \right)^{1/2} = (z^*z)^{1/2} = \| z \|_2$

↳ does not change length
or angles $(Qz)^*(Qy) = z^*y$

# Householder Transformation (reflection)

$$(1) \quad z = \begin{pmatrix} z_1 \\ \vdots \\ z_n \end{pmatrix} \in \mathbb{C}^n \qquad Hz = \begin{pmatrix} \gamma \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

where $|\gamma| = \|z\|_2$

$$H = I - 2\frac{vv^*}{v^*v} \qquad \text{or} \qquad H = I - 2vv^*$$

with $\|v\|_2 = 1$

$$H^{**}H = H^2 = \left(I - 2\frac{vv^*}{v^*v}\right)\left(I - 2\frac{vv^*}{v^*v}\right) =$$

$$I - 4\frac{vv^*}{v^*v} + 4\frac{vv^*vv^*}{(v^*v)^2} = I$$

Since $H = H^*$, obviously $HH^* = I$

$\to H$ unitary (orthog. if real)

How to pick $v$ for (1)   (fix $\gamma$ later)

(assume $\|v\| = 1$)

$$Hz = z - 2vv^*z = \gamma e_1 \implies$$
$$v(2 \cdot v^*z) = z - \gamma e_1$$

want $v$ normalized so scalar $2v^*z$ doesn't matter (doesn't change direction)

pick $v = z - \gamma e_1$ and normalize

(assuming $z \neq \gamma e_1 \to v = 0$)

Pick $\gamma$ to avoid poss. cancellation

$$\gamma = -\text{sign}(z_1)$$

## QR decomp.

$$H_1^* A = \begin{pmatrix} \alpha_1 & \tilde{a}_{12} & \tilde{a}_{13} & \cdots \\ 0 & & & \\ 0 & \tilde{a}_{22} & \vdots & \\ \vdots & & \vdots & \\ 0 & 1 & \vdots & \end{pmatrix} \qquad \tilde{a}_{ij} \rightarrow \text{changed}$$

$$H_2 \text{ s.t. } H_2(H_1 A) = \begin{pmatrix} \alpha_1 & \tilde{a}_{12} & * & \rule[0.5ex]{1cm}{0.4pt} & * \\ 0 & \alpha_2 & & & \\ \vdots & 0 & & & \\ \vdots & \vdots & & & \\ 0 & 0 & * & \rule[0.5ex]{1cm}{0.4pt} & * \end{pmatrix}$$

$\downarrow$

"*told touch only row 2 and below*

$$\hat{H}_2^{(n-1)\times(n-1)} \text{ s.t. } \hat{H}_2 \begin{pmatrix} \tilde{a}_{22} \\ \vdots \\ \vdots \\ \check{a}_{n2} \end{pmatrix} = \begin{pmatrix} \alpha_2 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

(as before)

$$H_2 = \left( \begin{array}{c|c} 1 & \\ \hline & \hat{H}_2 \end{array} \right) \qquad \underline{etc}$$

$$H_i \in \mathbb{C}^{n\times n} \quad H_n H_{n-1} H_{n-2} \cdots H_1 A = \begin{pmatrix} \alpha_1 & & \\ & \ddots & \\ 0 & & \alpha_m \end{pmatrix}$$
$$= R$$

$$A = H_1 H_2 \cdots H_{n-1} H_n R = QR$$

$$A^{n\times m} = Q^{n\times n} R^{n\times m} = Q^{n\times n} \begin{bmatrix} \begin{array}{c} R_m \\ \triangledown \\ 0 \end{array} \end{bmatrix}$$

$$= \begin{bmatrix} Q_1 & Q_2 \end{bmatrix}^{n\times m \; n\times(n-m)} \begin{bmatrix} R_m \\ 0 \end{bmatrix} = Q_1 R_m = Q_m R_m$$

$$r(Q_1) = r(A) \qquad r(Q_2) = r(A)^\perp = \text{null}(A^*)$$

$$\| b - Ax \|_2 = \| Q Q^* b - Q^* Q^* A x \|_2$$

$$= \left\| Q \left( Q^* b - \begin{bmatrix} R_m \\ 0 \end{bmatrix} x \right) \right\|_2 =$$

$$\left\| \begin{pmatrix} Q_1^* b \\ Q_2^* b \end{pmatrix} - \begin{pmatrix} R_m \\ 0 \end{pmatrix} x \right\|_2 \qquad \varphi$$

Solve $\quad Q_1^* b = R_m x$

$$\| b - Ax \|_2 = \| Q_2^* b \|_2$$

$$r = Q_2 Q_2^* b \qquad Q_1 Q_1^* b = Ax$$

---

cols of A not indep?

$$AP = \overset{\displaystyle m_1 \quad m_2 \quad n-m}{\left[ Q_1 \; Q_2 \; Q_3 \right]} \begin{bmatrix} R_1 & R_{02} \\ \diagdown & \square \\ 0 & 0 \\ 0 & 0 \end{bmatrix} \begin{matrix} m_1 \\ m_2 \\ m_3 \; n-m \end{matrix}$$

$$m = m_1 + m_2$$

avoiding overflow in computing $\alpha$

~~also~~ $\quad S = \dfrac{b}{(|a|^2+|b|^2)^{1/2}} = \dfrac{b/|b|}{\left(\dfrac{|a|^2}{|b|^2}+1\right)^{1/2}}$

$\quad = \dfrac{b/|a|}{\left(1+\dfrac{|b|^2}{|a|^2}\right)^{1/2}|a|} \quad \xleftarrow{\;\;} \quad \uparrow \; |b| \gtrless |a|$

$\qquad\qquad\qquad\qquad\qquad\quad |a| > |b|$

Same for computing $c$

$\tau = \dfrac{|a|}{|b|} \quad \rightarrow \quad S = \dfrac{b/|b|}{(\tau^2+1)^{1/2}}$

or

$\tau = \dfrac{|b|}{|a|} \quad \rightarrow \quad S = \dfrac{b/|a|}{(1+\tau^2)^{1/2}}$

same for $c$

Can use Givens rotations to selectively
set coeff. in vector (or matrix) to zero

$$\begin{pmatrix} 1 & & & & \\ & \ddots & & & \\ & & c & \bar{s} & \\ & & & 1 & \\ & & & & \ddots \\ & & -s & \bar{c} & \\ & & & & 1 \end{pmatrix} \begin{pmatrix} \{ \\ a_j \\ \} \\ a_j \\ \} \end{pmatrix} = \begin{pmatrix} \{ \\ \alpha \\ \} \\ 0 \\ \} \end{pmatrix}$$

# Givens Transformations (rotations)

real $\quad G = \begin{pmatrix} c & s \\ -s & c \end{pmatrix}$

complex $\quad G = \begin{pmatrix} c & \bar{s} \\ -s & \bar{c} \end{pmatrix}$ or $\begin{pmatrix} c & \bar{s} \\ -s & c \end{pmatrix}$ if

we pick $c \in \mathbb{R}$ (below)

verify $\begin{cases} G^*G = GG^* = I & \text{(complex)} \\ G^TG = GG^T = I & \text{(real)} \end{cases}$

$\begin{pmatrix} c & \bar{s} \\ -s & \bar{c} \end{pmatrix}\begin{pmatrix} a \\ b \end{pmatrix} = \begin{pmatrix} \alpha \\ 0 \end{pmatrix}$ $\quad$ where $\alpha = \left\| \begin{pmatrix} a \\ b \end{pmatrix} \right\|_2$

$\left( \text{or } |\alpha| = \| - \|_2 \right)$

$\begin{cases} ca + \bar{s}b = \alpha \\ -sa + \bar{c}b = 0 \end{cases}$ $\quad \left( \text{or } \alpha = \left\| \begin{pmatrix} a \\ b \end{pmatrix} \right\|_2 \frac{a}{|a|} \right)$

~~XXXXXXX~~ $\rightarrow$ if $a=0 \rightarrow \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$

$b=0 \rightarrow \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$

$\begin{cases} sca + |s|^2 b = s\alpha \\ -sca + |c|^2 b = 0 \end{cases}$ $\rightarrow b = s\alpha \implies s = b/\alpha$

(avoids underflow & over flow)

$-\frac{ab}{\alpha} + \bar{c}b = 0 \implies \bar{c} = \frac{a}{\alpha} \rightarrow c = \frac{\bar{a}}{\alpha}$

taking $\alpha = \left\| \begin{pmatrix} a \\ b \end{pmatrix} \right\|_2 \frac{a}{|a|} \rightarrow \bar{c} = \bar{\alpha} \cdot \frac{|a|}{\alpha} \left\| \begin{pmatrix} a \\ b \end{pmatrix} \right\|_2^{-1} \in \mathbb{R}$

$$A^{n \times m} = Q^{n \times n} R^{n \times m} = \hat{Q}^{n \times m} \hat{R}^{m \times m}$$

$$Q^* Q = Q Q^* = I \qquad R = \left(\begin{array}{c} \triangledown \\ O \end{array}\right) \begin{array}{l} ]m \\ ]n-m \end{array}$$

$$\hat{Q}^{n \times m} \hat{R}^{m \times m} \left\{ \begin{array}{l} \text{economy } QR \\ \text{thin } QR \\ \dots \end{array} \right.$$

If columns of A not indep $\rightarrow A = \hat{Q}^{n \times \tilde{m}} \hat{R}^{\tilde{m} \times m}$

$\rightarrow QR$ decomp w. pivoting (in general robust but not fail-safe)

$$\| Q^* A x - Q^* b \| = \| A x - b \|_2$$

$$\left( \| Q^* x \| = \left( x^* Q Q^* x \right)^{1/2} = \| x \| \quad \text{any } x \right.$$

$$\left\| \left(\begin{array}{c} \triangledown \\ O \end{array}\right) x = \left(\begin{array}{c} \overset{B_1}{b_1} \\ \overset{B_2}{b_2} \end{array}\right) \right\| \qquad \text{if col. A indep } \rightarrow R \text{ nonsing}$$

solve $R x = b_1$ $B_1$

no thing to be done about $b_2$ $B_2^2$ part

$$\min \| A x - b \|_2 = \min \left\| \left(\begin{array}{c} \triangledown \\ O \end{array}\right) x - \left(\begin{array}{c} b_1 \\ b_2 \end{array}\right) \right\|_2 = \| b_2 \|_2$$

$b_1 \in R(A) \rightarrow A x = b_1 \qquad b_2 \perp R(A)$

$$Q = \left[ \hat{Q}_1 \ Q_2 \right] \qquad \underset{n \times m}{A} = \underset{m \times m}{Q_1 R_1} \rightarrow \left(\begin{array}{c} Q_1^* \\ Q_2^* \end{array}\right) Q_1 R_1 =$$

$$Q_1 = \hat{Q} \qquad \left(\begin{array}{c} I \\ O \end{array}\right) R_1 \quad \overset{R_1}{\cancel{R_1}} \underset{O}{\cancel{O}} = \left(\begin{array}{c} R_1 \\ O \end{array}\right)$$

$$Q_1 Q_1^* = A(A^*A)^{-1}A^*$$

$$\rightarrow Q_1 R_1 \left( R_1^* Q_1^* Q_1 R_1 \right)^{-1} R_1^* Q_1^*$$

$$= Q_1^* R_1 \left( R_1^* R_1 \right)^{-1} R_1^* Q_1^*$$

$$= Q_1^* R_1 R_1^{-1} R_1^{-*} R_1^* Q_1^* = Q_1 Q_1^*$$

$$R(Q_1) = R(A)$$

$$R(Q_2) = R(A)^{\perp}$$

# Eigenvalues and Eigenvectors

Let $Ax = \lambda x$ and $y^* A = \lambda y^*$ (for same $\lambda$).

We call the vector $x$ a (right) eigenvector, the vector $y$ a left eigenvector, and $\lambda$ an eigenvalue of $A$, the triple together is called an eigentriple (of $A$), and $(\lambda, x)$ and $(\lambda, y)$ a (right) eigenpair and left eigenpair.

The set of all eigenvalues of $A$, $\Lambda(A)$, is called the spectrum of $A$ (when convenient we will count multiplicities in $\Lambda(A)$).

If the matrix $A$ is diagonalizable (has a complete set of eigenvectors) we have
$$A = V \Lambda V^{-1} \Leftrightarrow AV = V\Lambda,$$
where $V$ is a matrix with the right eigenvectors as columns and $\Lambda$ is a diagonal matrix with the eigenvalues as coefficients. This is not always possible (soon).

A similar decomposition can be given for the left eigenvectors.

# Spectral Radius

The spectral radius $\rho(A)$ is defined as $\rho(A) = \max\{|\lambda| : \lambda \in \Lambda(A)\}$.

Theorem:

For all $A$ and $\varepsilon > 0$ a consistent norm $\|\cdot\|_\alpha$ exists such that $\|A\|_\alpha \leq \rho(A) + \varepsilon$.

So, if $\rho(A) < 1$, then a consistent norm $\|\cdot\|_\alpha$ exists such that $\|A\|_\alpha < 1$.

Take $\varepsilon = \frac{1}{2}\left(1 - \rho(A)\right)$ and apply theorem above.

Define $A^* = \overline{A^T}$ (complex conjugate transpose).

If $A$ is Hermitian ($A = A^*$), then $\rho(A) = \|A\|_2$.

If $A$ is normal ($AA^* = A^*A$), then $\rho(A) = \|A\|_2$.

# Characteristic Polynomial

The eigenvalues of A are determined by the characteristic polynomial of A.

$$Ax = \lambda x \Longleftrightarrow (A - \lambda I)x = 0$$

So we're looking for (eigen)values $\lambda$ such that the matrix $(A - \lambda I)$ is singular:

$$\det(A - \lambda I) = 0 \qquad \text{(this is a polynomial in } \lambda)$$

This polynomial is called the characteristic polynomial of $A$. The eigenvalues of $A$ are defined to be the roots of its characteristic polynomial.

Since eigenvalues of matrix are roots of its characteristic polynomial, the Fundamental Theorem of Algebra implies that an $n \times n$ matrix $A$ always has $n$ eigenvalues. The eigenvalues, however, need be neither distinct nor real.

Complex eigenvalues of a real matrix must come in complex conjugate pairs.

# Multiplicity of eigenvalues

Eigenvalues may be single or multiple (single or multiple roots).
An eigenvalue with multiplicity $k > 1$ has k or fewer independent eigenvectors associated with it. It has at least one associated eigenvector. If it has fewer than k independent eigenvectors we call the eigenvalue (and the matrix) defective.

The multiplicity of an eigenvalue as the (multiple) root of the char. polynomial is called its algebraic multiplicity.
The number of independent eigenvectors associated with an eigenvalue is called its geometric multiplicity.

The geometric multiplicity is smaller than or equal to the algebraic multiplicity.

A matrix that is not defective is called diagonalizable: we have the decomposition

$$A = X \Lambda X^{-1} \Leftrightarrow X^{-1} A X = \Lambda = \mathrm{diag}(\lambda_i)$$

where $X$ contains the eigenvectors (as columns) and $\Lambda$ contains the eigenvalues.

# Jordan form of a matrix

For every matrix $A \in \mathbb{C}^{n \times n}$ there exists a nonsingular matrix $X$ such that

$$X^{-1}AX = \text{diag}(J_1, \cdots, J_q)$$

$$J_i = \begin{bmatrix} \lambda_i & 1 & & \\ & \lambda_i & \ddots & \\ & & \ddots & 1 \\ & & & \lambda_i \end{bmatrix} \text{ and } J_i \in \mathbb{C}^{m_i \times m_i}, \text{ and } m_1 + m_2 + \cdots + m_q = n.$$

Each block has one corresponding eigenvector: $q$ independent eigenvectors
Each block has $m_i - 1$ principal vectors (of grade 2)

If every block is of size 1, the matrix is diagonalizable
Multiple blocks can have the same eigenvalue: $\lambda_i = \lambda_j$
The sum of the sizes of all blocks with the same eigenvalue $\lambda$ is the algebraic multiplicity of the eigenvalue $\lambda$. The number of blocks with the same eigenvalue $\lambda$ is the geometric multiplicity of the eigenvalue $\lambda$.

# Invariant Subspaces

A generalization of eigenvectors to higher dimensions is an invariant subspace.

We say a subspace $\mathcal{V} \subset \mathbb{C}^n$ is invariant under $A \in \mathbb{C}^{n \times n}$ if

$$\text{for all } x \in V : Ax \in V$$

It is possible that an invariant subspace is the span of a set of eigenvectors, but this need not be the case. Moreover, in general, elements (vectors) of the subspace will not be themselves eigenvectors.

In many cases it is useful to consider the restriction of the matrix $A$ to the invariant subspace $\mathcal{V}$: $A_\mathcal{V} : \mathcal{V} \to \mathcal{V}$.

If $V \in \mathbb{C}^{n \times k}$ and $\text{range}(V) = \mathcal{V}$ then $AV = VL$ with $L \in \mathbb{C}^{k \times k}$.

Hence $L$ represents $A$ in the basis defined by $V$ for the space $\mathcal{V}$

# Similarity Transformation and Schur Decomposition

Let $A$ have eigenpairs $\left(\lambda_i, v_i\right)$: $Av_i = \lambda_i v_i$

For nonsingular $B$, define the *similarity transformation*: $BAB^{-1}$

The matrix $BAB^{-1}$ has the same eigenvalues, $\lambda_i$, as $A$ and eigenvectors $Bv_i$:

$$BAB^{-1}\left(Bv_i\right) = BAv_i = \lambda_i Bv_i$$

In fact, $BAB^{-1}$ has the same Jordan-block structure as $A$.

In many cases, we are interested in a (complex) unitary (or real, orthogonal) similarity transformation:

$$QAQ^* \qquad \text{with} \qquad QQ^* = Q^*Q = I$$

Schur decomposition: $QAQ^* = U$ (upper triangular)

# Similarity Transformation

Similarity transformation for $A$: $BAB^{-1}$;
this can be done with any nonsingular $B$.

Let $Ax = \lambda x$, then $BAB^{-1}Bx = BAx = \lambda Bx$.

$BAB^{-1}$ has the same eigenvalues as $A$, and
eigenvectors $Bx$ where $x$ is an eigenvector of $A$.

Although any nonsingular $B$ possible, most stable and
accurate algorithms with orthogonal (unitary) matrix.

For example used in the QR algorithm.

# Similarity Transformation

Orthogonal similarity transformation for $A$: $Q^*AQ$, where $Q^*Q = I$.

If $Q^*AQ = \begin{pmatrix} L_1 & F \\ 0 & L_2 \end{pmatrix}$, then $\mathcal{L}(A) = \mathcal{L}(L_1) \cup \mathcal{L}(L_2)$.

If we can find $Q \equiv \begin{bmatrix} Q_1 & Q_2 \end{bmatrix}$ that yields such a decomposition we have reduced the problem to two smaller problems.

Moreover, $AQ_1 = Q_1 L_1$ and $\mathrm{range}(Q_1)$ is invariant subspace. Eigenpair $L_1 z = \lambda z$ gives eigenpair $AQ_1 z = Q_1 L_1 z = \lambda Q_1 z$.

# Approximation over Search Space

For large matrices we cannot use full transformations. Often we do not need all eigenvalues/vectors. Look for proper basis $Q_1$ that captures relevant eigenpairs. We do not need $Q_2$.

Approximations over subspace $\mathrm{range}(Q_1)$: $L_1 = Q_1^* A Q_1$

When is an approximation good (enough)?

We will rarely find $A Q_1 - Q_1 L_1 = 0$ unless we do huge amount of work.
Not necessary. We are working with approximations and we must deal with numerical error anyway.

# Approximation over Search Space

Let $AQ_1 - Q_1L_1 = R$ with $\|R\|$ small relative to $\|A\|$.

Now, $\left(A - RQ_1^*\right)Q_1 - Q_1L_1 = AQ_1 - R - Q_1L_1 = 0$.

$\mathrm{range}\left(Q_1\right)$ is exact invariant subspace of perturbed matrix, $\hat{A}$.
$\hat{A} = A - RQ_1^*$ and $\left\|\hat{A} - A\right\|\big/\|A\| = \|R\|\big/\|A\|$

If $\|R\|\big/\|A\|$ sufficiently small, then $Q_1$ is acceptable.
In fact, this is as good as we can expect (unless we're lucky).
Any numerical operation involves perturbed operands!
Note that we cannot always say that $Q_1$ is accurate.

# Eigenvalue problems

Before we compute eigenvalues and eigenvectors numerically, we must understand what we can and cannot compute (accurately) or should not compute.

We may want $Ax = \lambda x$ for single eigenpair, for example with $\lambda$ as small as possible. Say minimum energy level.

In many cases we want $Ax_i = \lambda_i x_i$ for $i = 1 \ldots M$ ($M$ large), where $\lambda_i$ are smallest $M$ eigenvalues.

- It may be important that we do not skip any eigenvalues.
- We may want the invariant subspace accurately.
- We may want every eigenvector accurately.

# Usefulness of Computed Results

In general we need to consider the accuracy of a computed answer, without knowing the exact answer.

This involves the sensitivity of the result we want to compute.

If some result is very sensitive to small changes in the problem, it may be impossible to compute exactly. In other cases results may be computable but at very high price, for example, an algorithm may convergence very slowly.

Sometimes it is better to compute related but less sensitive result.

# Sensitivity of an Eigenvalue

Sensitivity of eigenvalues to perturbations in the matrix.

Different eigenvalues or eigenvectors of a matrix are not equally sensitive to perturbations of the matrix.

Let $Ax = \lambda x$ and $y^* A = \lambda y^*$, where $\|x\| = \|y\| = 1$.

Consider $(A + E)(x + e) = (\lambda + \varepsilon)(x + e)$ and drop second order terms.

$$Ax + Ae + Ex \cong \lambda x + \lambda e + \varepsilon x \Leftrightarrow Ae + Ex \cong \lambda e + \varepsilon x$$

$$y^* Ae + y^* Ex = \lambda y^* e + y^* Ex \cong \lambda y^* e + \varepsilon y^* x \quad \Leftrightarrow$$

$$y^* Ex \cong \varepsilon y^* x \Rightarrow \varepsilon \cong \frac{y^* Ex}{y^* x} \Rightarrow |\varepsilon| \leq \frac{\|E\|}{y^* x}$$

Condition number of simple eigenvalue: $\kappa(\lambda) = \left| y^* x \right|^{-1}$

# Sensitivity of an Eigenvalue

For symmetric/Hermitian matrix, right and left eigenvectors are the same. So, eigenvalues are inherently well-conditioned.

More generally, eigenvalues are well conditioned for normal matrices, but eigenvalues of nonnormal matrices need not be well conditioned.

Nonnormal matrices may not have a full set of eigenvectors. The algebraic multiplicity, the multiplicity of $\lambda$ as a root of $\det(A - \lambda I) = 0$, is not equal to the geometric multiplicity, $\dim \operatorname{null}(A - \lambda I)$. In that case we can consider conditioning of the invariant subspace associated with a Jordan block.

# Sensitivity of an Eigenvalue

If $\mu$ is an eigenvalue of $A + E$, then $\lambda \in \mathcal{L}(A)$ exists:

$$|\mu - \lambda| \leq \left\| XEX^{-1} \right\| \leq \kappa(X) \|E\|$$

where $X$ is the matrix of eigenvectors of $A$ and $\kappa(X) \equiv \|X\| \left\| X^{-1} \right\|$ is a condition number (consistent norm).

A useful backward error result is given by the residual.

Let $r = Ax - \lambda x$ and $\|x\| = 1$.

Then there exists a perturbation $E$ with $\|E\| \leq \|r\|$ such that

$$(A + E)x = \lambda x.$$

Proof: Take $E = -rx^*$.

# Sensitivity of Eigenvectors

Consider $A = \begin{pmatrix} 1 & 0 \\ 0 & 1+\varepsilon_1 \end{pmatrix}$ $\rightarrow$ $x_1 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$, $x_2 = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$.

Consider perturbation $E$ with $\|E\| = \varepsilon_1 + \varepsilon_2$, $\varepsilon_2$ arb. small.

Enough to give $A$ any eigenvectors not equal to $x_1$ and $x_2$.

Let $E = E_1 + E_2$ and $\hat{X} = \begin{bmatrix} \hat{x}_1 & \hat{x}_2 \end{bmatrix}$ (unitary)

Let $E_1 = \begin{pmatrix} 0 & 0 \\ 0 & -\varepsilon_1 \end{pmatrix}$ and $E_2 = \hat{X} \begin{pmatrix} 0 & 0 \\ 0 & \varepsilon_2 \end{pmatrix} \hat{X}^*$.

$A + E_1 = I$ (all nonzero vectors are eigenvectors)

$A + E = I + \hat{X} \begin{pmatrix} 0 & 0 \\ 0 & \varepsilon_2 \end{pmatrix} \hat{X}^* = \hat{X} \begin{pmatrix} 1 & 0 \\ 0 & 1+\varepsilon_2 \end{pmatrix} \hat{X}^*$

## SVD

$$A^{n \times m} = U \Sigma V^*$$

$U^{n \times n}$ unitary

$V^{m \times m}$ unitary

$\Sigma^{n \times m}$ diag, real, nonneg.

$$= U \begin{pmatrix} \sigma_1 & & \\ & \ddots & \sigma_m \\ 0 & & \end{pmatrix} V^* \quad \text{or} \quad U \begin{pmatrix} \sigma_1 & & \\ & \ddots & & 0 \\ & & \sigma_m \end{pmatrix} V^*$$

$$\text{or} \quad U \begin{pmatrix} \sigma_1 & \\ & \ddots \\ & & \sigma_m \end{pmatrix} V^* \quad \text{(square)}$$

$$\sigma_i \in \mathbb{R} \geq 0$$

convention: $\sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_m$

sometimes $\sigma_1 \geq \cdots \geq \sigma_p \qquad p = \min(m, n)$

or better $p = \text{rank}(A)$

$$A = U^{n \times m} \Sigma^{m \times m} (V^{m \times m})^* \qquad \text{economy SVD} \qquad (m \geq n)$$

$$\text{or} \quad U^{n \times n} \Sigma^{n \times n} V^{*\, n \times m}$$

If $A$ real $U, V$ real orthog.

$$R(U_p) = R(A) \qquad A = U^{n \times p} \Sigma^{p \times p} (V^*)^{p \times m}$$

$R(A) = u_1 \cdots u_p$      $\text{null}(A) = \text{span}\{v_{p+1}, v_m\}$

$R(A)^{\perp} = u_{p+1} \cdots u_n$      $\text{null}(A)^{\perp} = v_1 \cdots v_p$

$$A = \sum \sigma_i u_i v_i^* \qquad A^+ = \sum_{\substack{i: \\ \sigma_i > 0}} \sigma_i^{-1} v_i u_i^* = V \Sigma^+ U^*$$

$$x = A^{-1}b = \sum_{i=1}^{n} v_i \frac{u_i^* b}{\sigma_i}$$

(note sensitivity in $u_i$ directions corr. to smallest $\sigma_i$)

$$\|A\|_2 \|A^{-1}\|_2 = \sigma_1/\sigma_n \quad \text{or} \quad \sigma_{max}/\sigma_{min}$$

$$\|A\|_2 = \sigma_1 \quad (\text{max } \sigma_i)$$

$$\|A^{-1}\|_2 = \max\left(\frac{1}{\sigma_i}\right) = \sigma_{min}^{-1} = \sigma_n^{-1}$$

$$\|A\|_F = \left(\sum_i \sigma_i^2\right)^{1/2}$$

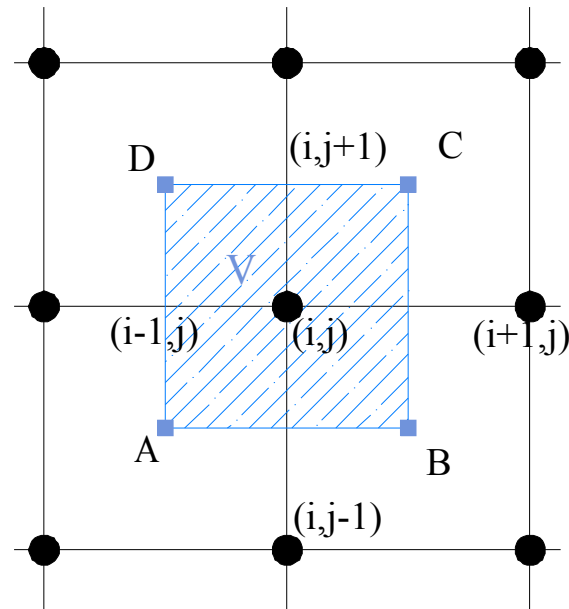$$\min_{A+E \text{ sing}} \|E\|_2 = \sigma_{min} \longrightarrow \text{distance to singularity}$$

---

with some effort (see notes corr. Watkins)

$$\frac{\|\Delta x\|_2}{\|x\|} \leq \left(\text{cond}(A)_2^2 \tan\theta + \text{cond}_2(A)\right)\frac{\|E\|_2}{\|A\|_2}$$

where $\theta \quad \angle(b, \text{range}(A))$

# Model Problems

Discretize $-\left(pu_x\right)_x - \left(qu_y\right)_y + ru_x + su_y + tu = f$.



Integrate equality over box $V$. Use Gauss' divergence theorem to get

$$\int_V \left(pu_x\right)_x + \left(qu_y\right)_y \, dx\, dy = \int_{\partial V} \begin{pmatrix} pu_x \\ qu_y \end{pmatrix} \cdot n \, ds$$

And approximate the line integral numerically.

# Model Problems

Now we approximate the boundary integral $\int_{\partial V} \begin{pmatrix} pu_x \\ qu_y \end{pmatrix} \cdot n \, ds$.

We approximate the integrals over each side of box $V$ using the midpoint rule and we approximate the derivatives using central differences.

$$\int_B^C pu_x n_1 \, dy \approx \frac{\Delta y}{\Delta x} p_{i+1/2,j} \left( U_{i+1,j} - U_{i,j} \right) \text{ and so on for the other sides}$$

We approximate the integrals over $ru_x$, $su_y$, $tu$, and $f$ using the area of the box and the value at the midpoint of the box, where we use central differences for derivatives. So, $u_x \approx \left( U_{i+1,j} - U_{i-1,j} \right) / \left( 2\Delta x \right)$, and so on.

For various examples we will also do this while strong convection relative to the mesh size makes central differences a poor choice (as it gives interesting systems).

**16**

# Model problems

This gives the discrete equations

$$-\frac{\Delta y}{\Delta x}\left[p_{i+1/2,j}\left(U_{i+1,j}-U_{i,j}\right)-p_{i-1/2,j}\left(U_{i,j}-U_{i-1,j}\right)\right]$$
$$-\frac{\Delta x}{\Delta y}\left[q_{i,j+1/2}\left(U_{i,j+1}-U_{i,j}\right)-p_{i,j-1/2}\left(U_{i,j}-U_{i,j-1}\right)\right]$$
$$+\left(\Delta y\,/\,2\right)r_{i,j}\left(U_{i+1,j}-U_{i-1,j}\right)+\left(\Delta x\,/\,2\right)s_{i,j}\left(U_{i,j+1}-U_{i,j-1}\right)$$
$$+\Delta x\Delta y\,t_{i,j}U_{i,j}=\Delta x\Delta y f_{i,j}$$

Often we divide this result again by $\Delta x\Delta y$.

# Rate of Convergence

Let $\hat{x}$ be the solution of $Ax = b$, and we have iterates $x_0, x_1, x_2, \ldots$

$\{x_k\}$ converges (q-)linearly to $\hat{x}$ if there are $N \geq 0$ and $c \in [0, 1)$ such that for $k \geq N : \| x_{k+1} - \hat{x} \| \leq c \| x_k - \hat{x} \|$,

$\{x_k\}$ converges (q-)superlinearly to $\hat{x}$ if there are $N \geq 0$ and a sequence $\{c_k\}$ that converges to $0$ such that for $k \geq N : \| x_{k+1} - \hat{x} \| \leq c_k \| x_k - \hat{x} \|$

$\{x_k\}$ converges to $\hat{x}$ with (q-)order at least $p$ if there are $p > 1$, $c \geq 0$, and $N \geq 0$ such that $k \geq N : \| x_{k+1} - \hat{x} \| \leq c \| x_k - \hat{x} \|^p$ (quadratic if $p = 2$, cubic if $p = 3$, and so on)

$\{x_k\}$ converges to $\hat{x}$ with j-step (q-)order at least $p$ if there are a fixed integer $j \geq 1$, $p > 1$, $c \geq 0$, and $N \geq 0$, such that $k \geq N : \| x_{k+j} - \hat{x} \| \leq c \| x_k - \hat{x} \|^p$