

1. Math

$\pi \approx 3.141\,59$ $e \approx 2.718\,28$ $\sqrt{2} \approx 1.414$ $\sqrt{3} \approx 1.732$
Binome, Trinome
 $(a \pm b)^2 = a^2 \pm 2ab + b^2$ $a^2 - b^2 = (a - b)(a + b)$
 $(a \pm b)^3 = a^3 \pm 3a^2b + 3ab^2 \pm b^3$
 $(a + b + c)^2 = a^2 + b^2 + c^2 + 2ab + 2ac + 2bc$

Folgen und Reihen
 $\sum_{k=1}^n k = \frac{n(n+1)}{2}$ $\sum_{k=0}^n q^k = \frac{1-q^{n+1}}{1-q}$ $\sum_{n=0}^{\infty} \frac{z^n}{n!} = e^z$
 Arithmetische Summenformel Geometrische Summenformel Exponentialreihe

Mittelwerte (\sum von i bis N) (Median: Mitte einer geordneten Liste)
 $\bar{x}_{ar} = \frac{1}{N} \sum x_i \geq \bar{x}_{geo} = \sqrt[N]{\prod x_i} \geq \bar{x}_{hm} = \frac{N}{\sum \frac{1}{x_i}}$
 Arithmetisches Geometrisches Mittel Harmonisches

Ungleichungen: Bernoulli-Ungleichung: $(1+x)^n \geq 1+nx$
 $||x| - |y|| \leq |x \pm y| \leq |x| + |y|$ $|\underline{x}^T \cdot \underline{y}| \leq \|\underline{x}\| \cdot \|\underline{y}\|$
 Dreiecksungleichung Cauchy-Schwarz-Ungleichung

Mengen: De Morgan: $\overline{A \cap B} = \overline{A} \cup \overline{B}$ $\overline{A \cup B} = \overline{A} \cap \overline{B}$

1.1. Exp. und Log. $e^x := \lim_{n \rightarrow \infty} (1 + \frac{x}{n})^n$ $e \approx 2,71828$
 $a^x = e^{x \ln a}$ $\log_a x = \frac{\ln x}{\ln a}$ $\ln x \leq x - 1$
 $\ln(x^a) = a \ln(x)$ $\ln(\frac{x}{a}) = \ln x - \ln a$ $\log(1) = 0$

1.2. Matrizen $\underline{A} \in \mathbb{K}^{m \times n}$
 $\underline{A} = (a_{ij}) \in \mathbb{K}^{m \times n}$ hat m Zeilen (Index i) und n Spalten (Index j)
 $(\underline{A} + \underline{B})^T = \underline{A}^T + \underline{B}^T$ $(\underline{A} \cdot \underline{B})^T = \underline{B}^T \cdot \underline{A}^T$
 $(\underline{A}^T)^{-1} = (\underline{A}^{-1})^T$ $(\underline{A} \cdot \underline{B})^{-1} = \underline{B}^{-1} \underline{A}^{-1}$

$\dim \mathbb{K} = n = \text{rang } \underline{A} + \dim \ker \underline{A}$ $\text{rang } \underline{A} = \text{rang } \underline{A}^T$
1.2.1. Quadratische Matrizen $\underline{A} \in \mathbb{K}^{n \times n}$
 regulär/invertierbar/nicht-singulär $\Leftrightarrow \det(\underline{A}) \neq 0 \Leftrightarrow \text{rang } \underline{A} = n$
 singulär/nicht-invertierbar $\Leftrightarrow \det(\underline{A}) = 0 \Leftrightarrow \text{rang } \underline{A} \neq n$

orthogonal $\Leftrightarrow \underline{A}^T = \underline{A}^{-1} \Rightarrow \det(\underline{A}) = \pm 1$
 symmetrisch: $\underline{A} = \underline{A}^T$ schief-symmetrisch: $\underline{A} = -\underline{A}^T$

1.2.2. Determinante von $\underline{A} \in \mathbb{K}^{n \times n}$: $\det(\underline{A}) = |\underline{A}|$
 $\det \begin{bmatrix} \underline{A} & \underline{0} \\ \underline{C} & \underline{D} \end{bmatrix} = \det \begin{bmatrix} \underline{A} & \underline{B} \\ \underline{0} & \underline{D} \end{bmatrix} = \det(\underline{A}) \det(\underline{D})$
 $\det(\underline{A}) = \det(\underline{A}^T)$ $\det(\underline{A}^{-1}) = \det(\underline{A})^{-1}$
 $\det(\underline{A}\underline{B}) = \det(\underline{A}) \det(\underline{B}) = \det(\underline{B}) \det(\underline{A}) = \det(\underline{B}\underline{A})$
 Hat \underline{A} 2 linear abhäng. Zeilen/Spalten $\Rightarrow |\underline{A}| = 0$

1.2.3. Eigenwerte (EW) λ und Eigenvektoren (EV) \underline{v}
 $\underline{A}\underline{v} = \lambda \underline{v}$ $\det \underline{A} = \prod \lambda_i$ $\text{Sp } \underline{A} = \sum a_{ii} = \sum \lambda_i$

Eigenwerte: $\det(\underline{A} - \lambda \underline{1}) = 0$ Eigenvektoren: $\ker(\underline{A} - \lambda_i \underline{1}) = \underline{v}_i$
 EW von Dreieck/Diagonal Matrizen sind die Elem. der Hauptdiagonale.

1.2.4. Spezialfall 2×2 Matrix \underline{A}
 $\det(\underline{A}) = ad - bc$ $\begin{bmatrix} a & b \\ c & d \end{bmatrix}^{-1} = \frac{1}{\det \underline{A}} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix}$
 $\text{Sp}(\underline{A}) = a + d$

$\lambda_{1/2} = \frac{\text{Sp } \underline{A}}{2} \pm \sqrt{\left(\frac{\text{Sp } \underline{A}}{2}\right)^2 - \det \underline{A}}$

1.2.5. Differentiation
 $\frac{\partial \underline{x}^T \underline{y}}{\partial \underline{x}} = \frac{\partial \underline{y}^T \underline{x}}{\partial \underline{x}} = \underline{y}$ $\frac{\partial \underline{x}^T \underline{A} \underline{x}}{\partial \underline{x}} = (\underline{A} + \underline{A}^T) \underline{x}$
 $\frac{\partial \underline{x}^T \underline{A} \underline{y}}{\partial \underline{A}} = \underline{x} \underline{y}^T$ $\frac{\partial \det(\underline{B} \underline{A} \underline{C})}{\partial \underline{A}} = \det(\underline{B} \underline{A} \underline{C}) (\underline{A}^{-1})^T$

1.2.6. Ableitungsregeln ($\forall \lambda, \mu \in \mathbb{R}$)
 Linearität: $(\lambda f + \mu g)'(x) = \lambda f'(x) + \mu g'(x)$
 Produkt: $(f \cdot g)'(x) = f'(x)g(x) + f(x)g'(x)$
 Quotient: $\left(\frac{f}{g}\right)'(x) = \frac{g(x)f'(x) - f(x)g'(x)}{g(x)^2}$ $\left(\frac{\text{NAZ} - \text{ZAN}}{N^2}\right)$
 Kettenregel: $(f(g(x)))' = f'(g(x))g'(x)$

1.3. Integrale $\int e^x dx = e^x = (e^x)'$
 Partielle Integration: $\int u w' = u w - \int u' w$
 Substitution: $\int f(g(x))g'(x) dx = \int f(t) dt$

$F(x) - C$	$f(x)$	$f'(x)$
$\frac{1}{q+1} x^{q+1}$	x^q	$q x^{q-1}$
$\frac{2\sqrt{ax^3}}{3}$	\sqrt{ax}	$\frac{a}{2\sqrt{ax}}$
$x \ln(ax) - x$	$\ln(ax)$	$\frac{1}{x}$
$\frac{1}{a^2} e^{ax} (ax - 1)$	$x \cdot e^{ax}$	$e^{ax} (ax + 1)$
$-\cos(x)$	$\sin(x)$	$\cos(x)$
$\cosh(x)$	$\sinh(x)$	$\cosh(x)$
$-\ln \cos(x) $	$\tan(x)$	$\frac{1}{\cos^2(x)}$

$\int e^{at} \sin(bt) dt = e^{at} \frac{a \sin(bt) + b \cos(bt)}{a^2 + b^2}$
 $\int \frac{dt}{\sqrt{at+b}} = \frac{2\sqrt{at+b}}{a}$ $\int t^2 e^{at} dt = \frac{(a-1)^2 + 1}{a^3} e^{at}$
 $\int t e^{at} dt = \frac{at-1}{a^2} e^{at}$ $\int x e^{ax^2} dx = \frac{1}{2a} e^{ax^2}$

1.3.1. Volumen und Oberfläche von Rotationskörpern um x-Achse
 $V = \pi \int_a^b f(x)^2 dx$ $O = 2\pi \int_a^b f(x) \sqrt{1 + f'(x)^2} dx$

2. Floating-Point Arithmetics

2.1. IEEE 754
 Representation of real number $x \in \mathbb{R}$:

$x = (-1)^s (1 + m 2^{-t}) 2^{e-b}$

Single Precision 1.8.23
 Double Precision 1.11.52
 Largest finite number $(2 - 2^{-t}) 2^{1 - \frac{em}{2}}$
 Smallest non zero Number $2^{1 - \frac{em}{2} - 1}$

Accuracy:
 $\frac{|x 2^{-x_1} - x_1|}{|x_1|} \leq 2^{-t}$

Machine Epsilon ϵ_M
 $f(x) = x(1 + \epsilon)$ $f(x) = \frac{x}{1 + \epsilon}$
 with different ϵ .

Maximum error is half the precision given by
 Machine Epsilon: $\epsilon_M = 2^{-t-1}$.

Worst case error
 $\rho = r(1 + \theta_n)$ with $|\theta_n| \leq \frac{n \epsilon_M}{1 - \frac{n}{2} \epsilon_M}$.
 The number of operations is given by n .

A more general formulation is:
 $\rho = r \prod_{i=1}^n (1 + \epsilon_i)^{\eta_i}$ with $\eta_i \in \{-1; +1\}$ this results in:
 $\rho = r(1 + \theta_n)$ with $|\theta_n| \leq \gamma_n$ and $\gamma_n = \frac{n \epsilon_M}{1 - \frac{n}{2} \epsilon_M}$

2.2. Complex Floating Point Arithmetic
 The error in complex floating point operations can also be represented with:

$x \odot y = (x \circ y)(1 + \epsilon)$
 with $|\epsilon| \leq \alpha \epsilon_M$. The alpha coefficients are given in the table below:

complex floating-point operation	\oplus	\ominus	\otimes	\oslash
α	1	1	$2^{3/2}$	6

Table 1.1. Factors for precision of complex floating-point operations

3. Linear Algebra

3.1. SVD
 A matrix is unitary if: $\underline{U}^H \underline{U} = \underline{U} \underline{U}^H = \underline{I}$

The singular decomposition (SVD) of the matrix $\underline{A} \in \mathbb{C}^{m \times n}$ reads as:
 $\underline{A} = \underline{U} \underline{\Sigma} \underline{V}^H$ where:
 \underline{U} and \underline{V} are unitary and $\underline{\Sigma} \in \mathbb{C}^{m \times n}$ is diagonal.

$\underline{\Sigma}_0 = \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_{\min(m,n)})$
 $\underline{\Sigma} = \begin{cases} \underline{\Sigma}_0, & \text{for } m = n \\ [\underline{\Sigma}_0, \underline{0}], & \text{for } m < n \text{ (wide matrix)} \\ \begin{bmatrix} \underline{\Sigma}_0 \\ \underline{0} \end{bmatrix}, & \text{for } m > n \text{ (tall matrix)} \end{cases}$

The SVD can also be represented in the following form:

$\underline{A} = \sum_{i=1}^{\min(m,n)} \sigma_i \underline{u}_i \underline{v}_i^H$

with
 $\underline{U} = [\underline{u}_1, \underline{u}_2, \dots, \underline{u}_m] \in \mathbb{C}^{m \times m}$
 $\underline{V} = [\underline{v}_1, \underline{v}_2, \dots, \underline{v}_n] \in \mathbb{C}^{n \times n}$

The rank of $r = \text{rank}(\underline{A})$ is given by the number of **non zero** singular values $\sigma_i > 0$
 $\text{range}(\underline{A}) = \text{span}(\underline{u}_1, \dots, \underline{u}_r)$
 $\text{null}(\underline{A}) = \text{span}(\underline{v}_{r+1}, \dots, \underline{v}_n)$

Removing the singular values that are zero $\sigma_{r+1} = \dots = \sigma_n = 0$ the reduced SVD can be expressed with:

$\underline{A} = \underline{U}_{red} \underline{\Sigma}_{red} \underline{V}_{red}^H$
 with
 $\underline{\Sigma}_{red} = \text{diag}(\sigma_1, \dots, \sigma_r)$, $\underline{U}_{red} = [\underline{u}_1, \dots, \underline{u}_r]$ and
 $\underline{V}_{red} = [\underline{v}_1, \dots, \underline{v}_r]$

Gram Matrix

$\underline{B} = \underline{A} \underline{A}^H$ is Hermitian by definition. It's SVD reads as:

$\underline{B} = \underline{U}_{red} \underline{\Sigma}_{red} \underline{U}_{red}^H$
 A similar discussion is possible for $\underline{C} = \underline{A}^H \underline{A}$:

$\underline{C} = \underline{A}^H \underline{A} = \underline{V}_{red} \underline{\Sigma}_{red} \underline{V}_{red}^H$

Both decompositions are like the (reduced) EVD.
 When computing the Gram matrix $\underline{A}^H \underline{A}$, every element is equal to the inner product of the two corresponding columns of \underline{A} . Let \underline{a}_i denote the i -th column of \underline{A} . Then, the i, j -element of the Gram matrix can be written as

$$[\underline{A}^H \underline{A}]_{ij} = \underline{a}_i^H \underline{a}_j$$

Only the upper triangular part of $\underline{A}^H \underline{A}$ must be computed due to the Hermitian structure of $\underline{A}^H \underline{A}$. The lower triangular are simply the complex conjugate of the corresponding elements of the upper triangular.
Operation Count Gram Matrix: #FLOPS = $\mathcal{O}(mn^2)$

3.2. Projectors
 Projectors are idempotent matrices. $\underline{P} \in \mathbb{C}^{m \times m}$ is a projector if:

$$\underline{P}^2 = \underline{P}$$

Orthogonal projectors are additionally Hermitian. Therefore:

$$\underline{P}^2 = \underline{P} \text{ and } \underline{P}^H = \underline{P}$$

is a orthogonal projector.

The orthogonal complement of a projector \underline{P} is defined as:
 $\underline{S}_P^\perp = \{\underline{v} : (\underline{I} - \underline{P})\underline{v} = \underline{v}\}$

The eigenvalues of a projector are out of the set $\lambda_i \in \{1; 0\}$. And therefore:

$$\underline{P} = \underline{U} \text{diag}(1, \dots, 1, 0, \dots, 0) \underline{U}^H$$

with reads in the reduced form:

$$\underline{P} = \underline{U}_{red} \underline{U}_{red}^H = \sum_{i=1}^r \underline{u}_i \underline{u}_i^H$$

with $\underline{U}_{red} = [\underline{u}_1, \dots, \underline{u}_r]$. Furthermore it can be concluded that:

$$\text{rank}(\underline{P}) = \text{tr}(\underline{P}) = r$$

3.2.1. Projectors of vectors
 We can construct a projector of a vector \underline{a} according to:

$$\underline{P}_a = \frac{\underline{a} \underline{a}^H}{\|\underline{a}\|_2^2}$$

The projector onto the orthogonal complement is given by:

$$\underline{P}_a^\perp = \underline{I} - \underline{P}_a = \underline{I} - \frac{\underline{a} \underline{a}^H}{\|\underline{a}\|_2^2}$$

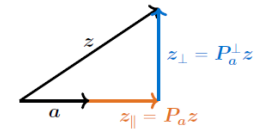


Fig. 2.1. Splitting a vector into a part along \underline{a} and perpendicular to \underline{a}

Given orthogonal vectors $\underline{a}_1, \dots, \underline{a}_n$ we can perform multiple projections onto the orthogonal complement via:

$$\prod_{i=1}^n \underline{P}_{\underline{a}_i}^\perp = \underline{I} - \sum_{i=1}^n \frac{\underline{a}_i \underline{a}_i^H}{\|\underline{a}_i\|_2^2}$$

3.3. Matrix and Vector norms
Vector norms:

p -norm: $\|\underline{x}\|_p = \sqrt[p]{\sum_{i=1}^n |x_i|^p}$
 1 -norm: $\|\underline{x}\|_1 = \sum_{i=1}^n |x_i|$
 2 -norm: $\|\underline{x}\|_2 = \sqrt{\sum_{i=1}^n |x_i|^2}$
 ∞ -norm: $\|\underline{x}\|_\infty = \max_{i \in \{1, \dots, n\}} |x_i|$

Matrix norms:
 $\|\underline{A}\| \leq \|\underline{A}\| \|\underline{B}\|$ The columnwise vectorization of \underline{A} reads as:
 $\text{vec}(\underline{A}) = [\underline{a}_1^T, \dots, \underline{a}_n^T]^T \in \mathbb{C}^{mn}$

$$\|\underline{A}\|_F = \|\text{vec}(\underline{A})\|_2 = \sqrt{\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2}$$

3.4. Householder Reflection

With the Householder Reflection we want remove all elements of a vector except the first one. Give a vector \mathbf{x} :

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \in \mathbb{C}^m$$

With the unitary matrix $\mathbf{F} \in \mathbb{C}^{m \times m}$ ($\mathbf{F}^H \mathbf{F} = \mathbf{F} \mathbf{F}^H = \mathbf{I}$) we want to achieve the following:

$$\mathbf{F} \mathbf{x} = \begin{bmatrix} \alpha \\ 0 \\ \vdots \\ 0 \end{bmatrix} = \alpha \mathbf{e}_1$$

It follows: $\|\mathbf{F} \mathbf{x}\|_2^2 = \mathbf{x}^H \mathbf{F}^H \mathbf{F} \mathbf{x} = \|\mathbf{x}\|_2^2$ and therefore: $|\alpha| = \|\mathbf{x}\|_2^2$

The overall goal is to find the matrix \mathbf{F} that achieves:

$$\mathbf{F} \mathbf{x} = \begin{bmatrix} \|\mathbf{x}\|_2 e^{j\phi} \\ 0 \\ \vdots \\ 0 \end{bmatrix} = \|\mathbf{x}\|_2 e^{j\phi} \mathbf{e}_1$$

We want to find the vector

$$\mathbf{v} = \|\mathbf{x}\|_2 e^{j\phi} \mathbf{e}_1 - \mathbf{x}$$

The vectors \mathbf{x} , \mathbf{v} and $\|\mathbf{x}\|_2 e^{j\phi} \mathbf{e}_1$ form an isosceles triangle. The height vector is given by the projection onto the orthogonal complement of \mathbf{v} :

$$\mathbf{h} = \mathbf{P}_{\mathbf{v}}^\perp \mathbf{x} = (\mathbf{I} - \mathbf{P}_{\mathbf{v}}) \mathbf{x} = \mathbf{x} - \frac{\mathbf{v} \mathbf{v}^H}{\|\mathbf{v}\|_2^2} \mathbf{x}$$

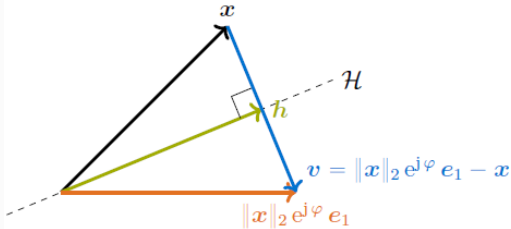


Fig. 2.2. Householder reflection

The matrix \mathbf{F} is given by:

$$\mathbf{F} = \mathbf{I} - 2 \frac{\mathbf{v} \mathbf{v}^H}{\|\mathbf{v}\|_2^2}$$

Decomposing the Householder reflector into it's EVD leads to:

$$[\mathbf{w}, \mathbf{W}_\perp] \text{diag}(-1, 1, \dots, 1) [\mathbf{w}, \mathbf{W}_\perp]^H$$

where $\mathbf{w} = \frac{\mathbf{v}}{\|\mathbf{v}\|_2}$ and $\mathbf{P}_{\mathbf{v}}^\perp = \mathbf{W}_\perp \text{erp} \mathbf{W}_\perp^H$. Only $\lambda_1 = -1$ of \mathbf{F} all other eigenvalues are $\lambda_i = 1$.

3.4.1. Choice of Phase after Reflection

A possible choice is: $\phi = 0$ such that α is positive-real. An other goal is to make \mathbf{v} as long as possible therefore:

$$\phi_{opt} = \arg \max_{\phi} \|\mathbf{v}\|_2^2$$

with $\mathbf{x} = [x_1, x_2]^T$ and

$\mathbf{v} = \|\mathbf{x}\|_2 e^{j\phi} \mathbf{e}_1 - \mathbf{x}$ the objective function can be expressed as

$$\|\mathbf{v}\|_2^2 = |e^{j\phi} \|\mathbf{x}\|_2 - x_1|^2 + \|x_2\|_2^2$$

As the second term is independent of ϕ , the first term must be maximized. To this end, it is necessary to choose $\phi_{opt} = \psi + \pi$ for $x_1 = |x_1| e^{j\psi}$, resulting in

$$\mathbf{v} = \begin{bmatrix} -(\|\mathbf{x}\|_2 + |x_1|) e^{j\psi} \\ -x_2 \end{bmatrix}$$

The maximum objective function is:

$$\|\mathbf{v}_{opt}\|_2^2 = 2 \|\mathbf{x}_2\|_2^2 + 2|x_1| \|\mathbf{x}\|_2$$

4. QR Decomposition

For a full-rank and tall matrix $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_n] \in \mathbb{C}^{m \times n}$ with $m \gg n$ the QR decomposition reads as:

$$\mathbf{A} = \mathbf{Q} \mathbf{R}$$

where $\mathbf{Q} \in \mathbb{C}^{m \times m}$ is unitary and

$$\mathbf{R} = \begin{bmatrix} \mathbf{R}_{red} \\ \mathbf{0} \end{bmatrix} \in \mathbb{C}^{m \times n}$$

with

$$\mathbf{R}_{red} = \begin{bmatrix} r_{11} & r_{12} & \dots & r_{1n} \\ 0 & r_{22} & \dots & r_{2n} \\ 0 & 0 & \ddots & \vdots \\ 0 & \dots & r_{nn} \end{bmatrix} \in \mathbb{C}^{n \times n}$$

with ($\nu \leq n$):

$$\text{span}(\mathbf{a}_1, \dots, \mathbf{a}_\nu) = \text{span}(\mathbf{q}_1, \dots, \mathbf{q}_\nu)$$

which results in:

$$\text{range}(\mathbf{A}) = \text{span}(\mathbf{a}_1, \dots, \mathbf{a}_n) = \text{span}(\mathbf{q}_1, \dots, \mathbf{q}_n) = \text{range}(\mathbf{Q}_{red})$$

\mathbf{Q}_{red} is an orthonormal basis for the range space of \mathbf{A} .

$$\mathbf{Q}_{red} = [\mathbf{q}_1, \dots, \mathbf{q}_n] \in \mathbb{C}^{m \times n}$$

$$\mathbf{Q}_{red}^\perp = [\mathbf{q}_{n+1}, \dots, \mathbf{q}_m] \in \mathbb{C}^{m \times m-n}$$

It follows from $\mathbf{Q} = [\mathbf{Q}_{red}, \mathbf{Q}_{red}^\perp]$ that

$$\mathbf{I} = \mathbf{Q}^H \mathbf{Q} =$$

$$\begin{bmatrix} \mathbf{Q}_{red}^H \\ \mathbf{Q}_{red}^{\perp H} \end{bmatrix} \cdot [\mathbf{Q}_{red}, \mathbf{Q}_{red}^\perp] = \begin{bmatrix} \mathbf{Q}_{red}^H \mathbf{Q}_{red} & \mathbf{Q}_{red}^H \mathbf{Q}_{red}^\perp \\ \mathbf{Q}_{red}^{\perp H} \mathbf{Q}_{red} & \mathbf{Q}_{red}^{\perp H} \mathbf{Q}_{red}^\perp \end{bmatrix}$$

Therefore: $\mathbf{Q}_{red}^H \mathbf{Q}_{red}^\perp = \mathbf{0}$. Consequently the space spanned by the columns of \mathbf{Q}_{red}^\perp is perpendicular to $\text{range}(\mathbf{Q}_{red})$. In other words:

$$\text{null}(\mathbf{A}^H) = \text{range}(\mathbf{Q}_{red}^\perp)$$

that is, the matrix \mathbf{Q}_{red}^\perp is a basis for the left nullspace of \mathbf{A} .

4.1. Operation Count of basic operations

f	Mult.+Sum.	#FLOPS
$\mathbf{a}^H \mathbf{b} = \sum_{i=1}^m a_i^* b_i$	$m + m - 1$	$2m - 1$
$\ \mathbf{a}\ _2 = \sqrt{\sum_{i=1}^m a_i^* a_i}$	$m + m - 1 + 1$	$2m$
$\gamma \mathbf{a}$	m	m
$\mathbf{a} \pm \mathbf{b}$	m	m

4.2. Classical Gram-Schmidt Orthogonalization

Algorithm 1 Classical Gram-Schmidt Orthogonalization
1: function CGS(A) \triangleright with $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_n] \in \mathbb{C}^{m \times n}$ with $m \geq n$
2: for $j = 1$ to n do
3: $\mathbf{v}_j \leftarrow \mathbf{a}_j$
4: for $i = 1$ to $j - 1$ do \triangleright orthogonalize the j -th column
5: $r_{ij} \leftarrow \mathbf{q}_i^H \mathbf{a}_j$ \triangleright coefficient of projection onto \mathbf{q}_i
6: $\mathbf{v}_j \leftarrow \mathbf{v}_j - r_{ij} \mathbf{q}_i$ \triangleright remove component along \mathbf{q}_i
7: end for
8: $r_{jj} \leftarrow \|\mathbf{v}_j\|_2$
9: $\mathbf{q}_j \leftarrow \frac{\mathbf{v}_j}{r_{jj}}$ \triangleright normalization
10: end for
11: return $\mathbf{Q}_{red}, \mathbf{R}_{red}$
12: end function

4.3. Modified Gram-Schmidt Orthogonalization

Algorithm 2 Modified Gram-Schmidt Orthogonalization
1: function MGS1(A) \triangleright with $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_n] \in \mathbb{C}^{m \times n}$ with $m \geq n$
2: for $j = 1$ to n do
3: $\mathbf{v}_j \leftarrow \mathbf{a}_j$
4: for $i = 1$ to $j - 1$ do \triangleright from CGS to MGS: change \mathbf{a}_j to \mathbf{v}_j
5: $r_{ij} \leftarrow \mathbf{q}_i^H \mathbf{v}_j$
6: $\mathbf{v}_j \leftarrow \mathbf{v}_j - r_{ij} \mathbf{q}_i$
7: end for
8: $r_{jj} \leftarrow \|\mathbf{v}_j\|_2$
9: $\mathbf{q}_j \leftarrow \frac{\mathbf{v}_j}{r_{jj}}$
10: end for
11: return $\mathbf{Q}_{red}, \mathbf{R}_{red}$
12: end function

Algorithm 3 Modified Gram-Schmidt Orthogonalization
1: function MGS2(A) $\triangleright \mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_n]$
2: for $j = 1$ to n do
3: $\mathbf{v}_j \leftarrow \mathbf{a}_j$ \triangleright store columns of \mathbf{A} in $\mathbf{v}_1, \dots, \mathbf{v}_n$
4: end for
5: for $i = 1$ to n do
6: $r_{ii} \leftarrow \|\mathbf{v}_i\|_2$
7: $\mathbf{q}_i \leftarrow \frac{\mathbf{v}_i}{r_{ii}}$ \triangleright normalization
8: for $j = i + 1$ to n do
9: $r_{ij} \leftarrow \mathbf{q}_i^H \mathbf{v}_j$ \triangleright coefficient of projection onto \mathbf{q}_i
10: $\mathbf{v}_j \leftarrow \mathbf{v}_j - r_{ij} \mathbf{q}_i$ \triangleright subtraction of component along \mathbf{q}_i
11: end for
12: end for
13: return $\mathbf{Q}_{red}, \mathbf{R}_{red}$
14: end function

Operation Count: #FLOPS = $\mathcal{O}(2mn^2)$

4.4. Householder QR Decomposition

With the Householder reflector $\mathbf{Q}_1 = \mathbf{F}_1$ the first column of matrix \mathbf{A} is transformed as follows

$$\mathbf{Q}_1 \mathbf{A} = \begin{bmatrix} \times & \times & \times \\ 0 & \times & \times \\ 0 & \times & \times \\ 0 & \times & \times \\ 0 & \times & \times \end{bmatrix}$$

In the next step we apply the matrix \mathbf{Q}_2 which inherits a Householder reflector $\mathbf{F}_2 \in \mathbb{C}^{m-1 \times m-1}$

$$\mathbf{Q}_2 = \begin{bmatrix} 1 & \mathbf{0}^T \\ \mathbf{0} & \mathbf{F}_2 \end{bmatrix} \in \mathbb{C}^{m \times m}$$

After that the matrix $\mathbf{Q}_2 \mathbf{Q}_1 \mathbf{A}$ looks like this

$$\mathbf{Q}_2 \mathbf{Q}_1 \mathbf{A} = \begin{bmatrix} \times & \times & \times \\ 0 & \times & \times \\ 0 & 0 & \times \\ 0 & 0 & \times \\ 0 & 0 & \times \end{bmatrix}$$

The unitary rotation matrix \mathbf{Q} at step $i \in [1, \dots, n]$ is given by

$$\mathbf{Q}_i = \begin{bmatrix} \mathbf{I}_{i-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{F}_i \end{bmatrix} \in \mathbb{C}^{m \times m}$$

with $\mathbf{F}_i = \mathbf{I} - 2 \frac{\mathbf{v}_i \mathbf{v}_i^H}{\|\mathbf{v}_i\|_2^2} \in \mathbb{C}^{m-i+1 \times m-i+1}$

After rotating all n columns of \mathbf{A} we generate the QR decomposition with $\mathbf{A} = \mathbf{Q}_1^H \mathbf{Q}_2^H \dots \mathbf{Q}_n^H \mathbf{R} = \mathbf{Q} \mathbf{R}$

$\mathbf{R} \in \mathbb{C}^{m \times n}$

Furthermore it holds

$\mathbf{Q} = \mathbf{Q}_1^H \mathbf{Q}_2^H \dots \mathbf{Q}_n^H = \mathbf{Q}_1 \mathbf{Q}_2 \dots \mathbf{Q}_n$ since \mathbf{Q}_i is Hermitian.

Algorithm 4 Householder QR Decomposition

1: function HOUSEHOLDERQR(A) $\triangleright \mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_n] \in \mathbb{C}^{m \times n}$ with $m \geq n$
2: for $i = 1$ to n do
3: $\mathbf{x} \leftarrow [\mathbf{a}]_{i:m,i}$ \triangleright take entries i to m of the i -th column
4: $[\mathbf{A}]_{i:m,i} \leftarrow \exp(j \arg(x_1) + \pi) |x|_2 \mathbf{e}_1$ \triangleright result of reflection for i -th column
5: $\mathbf{v}_i \leftarrow [\mathbf{A}]_{i:m,i} - \mathbf{x}$
6: $\mathbf{v}_i \leftarrow \frac{\mathbf{v}_i}{\|\mathbf{v}_i\|_2}$ \triangleright normalization
7: $[\mathbf{A}]_{i:m,i+1:n} \leftarrow [\mathbf{A}]_{i:m,i+1:n} - 2\mathbf{v}_i (\mathbf{v}_i^H [\mathbf{A}]_{i:m,i+1:n})$ \triangleright no reflection for $i = n$
8: end for
9: return $\mathbf{A}, \mathbf{v}_1, \dots, \mathbf{v}_n$
10: end function

Operation Count Algorithm 4: #FLOPS = $\mathcal{O}(2mn^2 - \frac{2}{3}n^3)$

Algorithm 5 Householder QR Decomposition: Computation of \mathbf{Q}_{red}

1: function HOUSEHOLDERQRQ($\mathbf{v}_1, \dots, \mathbf{v}_n$) $\triangleright \mathbf{v}_i$ defines the rotation \mathbf{Q}_i
2: $\mathbf{Q} \leftarrow \begin{bmatrix} \mathbf{I}_n \\ \mathbf{0} \end{bmatrix}$ \triangleright use first n columns of \mathbf{I} as initialization
3: for $i = n$ downto 1 do
4: $[\mathbf{Q}]_{i:m,i:n} \leftarrow [\mathbf{Q}]_{i:m,i:n} - 2\mathbf{v}_i (\mathbf{v}_i^H [\mathbf{Q}]_{i:m,i:n})$ \triangleright reflection based on \mathbf{v}_i
5: end for
6: return \mathbf{Q} \triangleright return \mathbf{Q}_{red}
7: end function

Algorithm 6 Householder QR Decomposition: Computation of $\mathbf{Q}^H \mathbf{b}$

1: function HOUSEHOLDERQRQ($\mathbf{v}_1, \dots, \mathbf{v}_n, \mathbf{b}$) $\triangleright \mathbf{v}_i$ defines the rotation \mathbf{Q}_i
2: for $i = 1$ to n do
3: $[\mathbf{b}]_{i:m} \leftarrow [\mathbf{b}]_{i:m} - 2\mathbf{v}_i (\mathbf{v}_i^H [\mathbf{b}]_{i:m})$ \triangleright reflection with normalized \mathbf{v}_i
4: end for
5: return \mathbf{b}
6: end function

Operation Count Algorithm 5: #FLOPS = $\mathcal{O}(2mn^2 - \frac{2}{3}n^3)$

Operation Count Algorithm 6: #FLOPS = $\mathcal{O}(4mn - 2n^2)$

Algorithm 7 Householder QR Decomposition: Computation of \mathbf{Q}_{red}^\perp

1: function HOUSEHOLDERQRPER($\mathbf{v}_1, \dots, \mathbf{v}_n$) $\triangleright \mathbf{v}_i$ defines the rotation \mathbf{Q}_i
2: $\mathbf{Q} \leftarrow \begin{bmatrix} \mathbf{0} \\ \mathbf{I}_{m-n} \end{bmatrix}$ \triangleright use last $m - n$ columns of \mathbf{I} as initialization
3: for $i = n$ downto 1 do
4: $[\mathbf{Q}]_{i:m,1:m-n} \leftarrow [\mathbf{Q}]_{i:m,1:m-n} - 2\mathbf{v}_i (\mathbf{v}_i^H [\mathbf{Q}]_{i:m,1:m-n})$ \triangleright reflection based on \mathbf{v}_i
5: end for
6: return \mathbf{Q} \triangleright return \mathbf{Q}_{red}^\perp
7: end function

Operation Count Algorithm 7: #FLOPS = $\mathcal{O}(4m^2n + 2n^3 - 6mn^2)$
If $m = n$ the number of FLOPS for Algorithm 7 is zero, because the partition \mathbf{Q}_{red}^\perp does not exist in this case.

4.4.1. Householder for null(\mathbf{A}^H)

In the case, that a basis for the nullspace \mathbf{A}^H , i.e., the orthogonal complement of $\text{range}(\mathbf{A})$ must be found Algorithm 7 can be employed. $\text{null}(\mathbf{A}^H) = \text{range}(\mathbf{Q}_{red}^\perp)$

$$\mathbf{Q} = [\mathbf{Q}_{red}, \mathbf{Q}_{red}^\perp]$$

5. Back Substitution

Solving an equation system with an upper triangular \mathbf{R} . $\mathbf{R}\mathbf{x} = \mathbf{b}$ can be solved with Algorithm 8.

$$\mathbf{R} = \begin{bmatrix} r_{11} & r_{12} & \cdots & r_{1n} \\ 0 & r_{22} & & \\ & & \ddots & \\ & & & r_{nn} \end{bmatrix} \in \mathbb{C}^{n \times n}$$

Algorithm 8 Back Substitution: Computation of the solution to $R\mathbf{x} = \mathbf{b}$

```

1: function BACKSUBSTITUTION( $R, \mathbf{b}$ )  $\triangleright R \in \mathbb{C}^{n \times n}$  is upper triangular
2:   for  $i = n$  downto 1 do
3:      $x_i \leftarrow \frac{1}{r_{ii}} \left( b_i - \sum_{j=i+1}^n r_{ij} x_j \right)$   $\triangleright$  for  $i = n$ , the summation is zero
4:   end for
5:   return  $\mathbf{x}$ 
6: end function
```

Operation Count Algorithm 8: #FLOPS = $\mathcal{O}(n^2)$

Algorithm 9 Back Substitution: Computation of R^{-1}

```

1: function BACKSUBSTITUTIONINVERSE( $R$ )  $\triangleright R \in \mathbb{C}^{n \times n}$  is upper triangular
2:    $P \leftarrow \mathbf{I}$   $\triangleright$  apply back substitution to columns of the identity matrix
3:   for  $j = 1$  to  $n$  do
4:     for  $i = j$  downto 1 do
5:        $p_{ij} \leftarrow \frac{1}{r_{ii}} \left( p_{ij} - \sum_{d=i+1}^j r_{id} p_{dj} \right)$   $\triangleright$  for  $i = n$ , the summation is zero
6:     end for
7:   end for
8:   return  $P$ 
9: end function
```

Operation Count Algorithm 9: #FLOPS = $\mathcal{O}(\frac{1}{3}n^3)$

$$\begin{aligned} \text{\#FLOPS} &= \sum_{j=1}^n \sum_{i=1}^j (2j - 2i + 1) \\ &= \sum_{j=1}^n [2j^2 - j(j+1) + j] = \frac{(2n+1)(n+1)n}{6} \end{aligned}$$

6. Least Squares

6.1. Least Square for Full-Rank Matrix

$\mathbf{A}\mathbf{x} = \mathbf{b}$ with $\mathbf{A} \in \mathbb{C}^{m \times n}$ and $m \geq n$. For tall $\mathbf{A} \in \mathbb{C}^{m \times n}$ with $m > n$ no solution exists in general, except $\mathbf{b} \in \text{range}(\mathbf{A})$.

For a square and invertible \mathbf{A} the solution

The Least Square solution is given by the normal equation:

$$\mathbf{A}^H \mathbf{A} \mathbf{x}_{\text{LS}} = \mathbf{A}^H \mathbf{b}$$

$$\mathbf{x}_{\text{LS}} = (\mathbf{A}^H \mathbf{A})^{-1} \mathbf{A}^H \mathbf{b}$$

The pseudoinverse of \mathbf{A} : $\mathbf{A}^\dagger = (\mathbf{A}^H \mathbf{A})^{-1} \mathbf{A}^H \mathbf{b}$

6.2. Least Square with Cholesky

We need to compute the Gram matrix: $\mathbf{A}^H \mathbf{A} = \mathbf{L} \mathbf{D} \mathbf{L}^H$.

$$\mathbf{L} \mathbf{D} \mathbf{L}^H \mathbf{x}_{\text{LS}} = \mathbf{A}^H \mathbf{b}$$

Due to \mathbf{L} being lower triangular we need two times backsubstitution.

The condition number with Cholesky increases because if we compute the Gram matrix it holds:

$$\kappa(\mathbf{A}^H \mathbf{A}) = \kappa^2(\mathbf{A})$$

The condition number of the Gram matrix of \mathbf{A} is the square of the original condition number κ .

6.3. Least Square with QR

Using the decomposition $\mathbf{A} = \mathbf{Q}_{\text{red}} \mathbf{R}_{\text{red}}$ We obtain the solution

$$\mathbf{A}^\dagger = \mathbf{R}_{\text{red}}^{-1} \mathbf{Q}_{\text{red}}^H$$

$$\mathbf{x}_{\text{LS}} = \mathbf{A}^\dagger \mathbf{b}$$

Where $\mathbf{R}_{\text{red}}^{-1}$ describes a backsubstitution. Therefore we need to compute $\mathbf{Q}^H \mathbf{b}$ and then employ a backsubstitution on this result. QR decomposition and Cholesky have the same order of complexity if $m = n$.

Algorithm 10 Least Squares via Cholesky Factorization

```

1: function LSCHOLESKY( $\mathbf{A}, \mathbf{b}$ )  $\triangleright$  constituting matrix  $\mathbf{A} \in \mathbb{C}^{m \times n}$  and right-hand side  $\mathbf{b}$ 
2:    $\mathbf{B} \leftarrow \mathbf{A}^H \mathbf{A}$   $\triangleright$  compute Gram of  $\mathbf{A}$ 
3:    $\mathbf{c} \leftarrow \mathbf{A}^H \mathbf{b}$   $\triangleright$  right-hand side of normal equations
4:    $\mathbf{L}, \mathbf{D} \leftarrow$  Cholesky factorization of  $\mathbf{B}$ 
5:    $\mathbf{v} \leftarrow \mathbf{L}^{-1} \mathbf{c}$   $\triangleright$  with back substitution
6:    $\mathbf{w} \leftarrow \mathbf{D}^{-1} \mathbf{v}$   $\triangleright$  scaling every element
7:    $\mathbf{x}_{\text{LS}} \leftarrow \mathbf{L}^{H-1} \mathbf{w}$   $\triangleright$  with back substitution
8: end function
```

Operation Count Algorithm 10: #FLOPS = $\mathcal{O}(mn^2 + \frac{1}{3}n^3)$
Computing the Gram matrix and Cholesky factorization is dominant. In total we have: $\mathcal{O}(mn^2 + 2mn + \frac{1}{3}n^3 + 2n^2 + n)$

Algorithm 11 Least Squares via QR Decomposition

```

1: function LSQR( $\mathbf{A}, \mathbf{b}$ )  $\triangleright$  constituting matrix  $\mathbf{A} \in \mathbb{C}^{m \times n}$  and right-hand side  $\mathbf{b}$ 
2:    $\mathbf{Q}, \mathbf{R} \leftarrow$  QR decomposition of  $\mathbf{A}$   $\triangleright$  preferably Householder QR
3:    $\mathbf{c} \leftarrow \mathbf{Q}^H \mathbf{b}$ 
4:    $\mathbf{x}_{\text{LS}} \leftarrow \mathbf{R}_{\text{red}}^{-1} \mathbf{c}$   $\triangleright$  with back substitution
5: end function
```

Operation Count Algorithm 11: #FLOPS = $\mathcal{O}(2mn^2 - \frac{2}{3}n^3)$
The complexity analysis assumes that we used Householder QR decomposition.

6.4. Least Squares with SVD

We decompose $\mathbf{A} \in \mathbb{C}^{m \times n}$ as follows

$$\mathbf{A} = \mathbf{U}_{\text{red}} \Sigma_{\text{red}} \mathbf{V}^H$$

with sub-unitary $\mathbf{U}_{\text{red}} \in \mathbb{C}^{m \times n}$, diagonal $\Sigma_{\text{red}} \in \mathbb{C}^{n \times n}$ and unitary $\mathbf{V} \in \mathbb{C}^{n \times n}$. The pseudoinverse of \mathbf{A} can be written as

$$\mathbf{A}^\dagger = \mathbf{V} \Sigma_{\text{red}}^{-1} \mathbf{U}_{\text{red}}^H$$

Thus the solution can be expressed as

$$\mathbf{x}_{\text{LS}} = \mathbf{V} \Sigma_{\text{red}}^{-1} \mathbf{U}_{\text{red}}^H \mathbf{b}$$

Algorithm 12 Least Squares via SVD

```

1: function LSSVD( $\mathbf{A}, \mathbf{b}$ )  $\triangleright$  constituting matrix  $\mathbf{A} \in \mathbb{C}^{m \times n}$  and right-hand side  $\mathbf{b}$ 
2:    $\mathbf{U}_{\text{red}}, \Sigma_{\text{red}}, \mathbf{V} \leftarrow$  singular value decomposition of  $\mathbf{A}$ 
3:    $\mathbf{c} \leftarrow \mathbf{U}_{\text{red}}^H \mathbf{b}$ 
4:    $\mathbf{w} \leftarrow \Sigma_{\text{red}}^{-1} \mathbf{c}$ 
5:    $\mathbf{x}_{\text{LS}} \leftarrow \mathbf{V} \mathbf{w}$ 
6: end function
```

Operation Count Algorithm 12: #FLOPS = $\mathcal{O}(2mn^2 + 11n^3)$

6.5. Comparison of Different Approaches

$$\text{\#FLOPS}_{\text{Cholesky}} \leq \text{\#FLOPS}_{\text{QR-Householder}} \leq \text{\#FLOPS}_{\text{SVD}}$$

The left inequality becomes an equality for $m = n$ and the right inequality is approximately an equality for $m \gg n$.

For numerical stability, the order is the same

$$\text{stability}_{\text{Cholesky}} \leq \text{stability}_{\text{QR-Householder}} \leq \text{stability}_{\text{SVD}}$$

Due to the **explicit computation of the Gram matrix** $\mathbf{A}^H \mathbf{A}$ in the Cholesky factorization based approach the condition of the problem is deteriorated considerably.

QR decomposition does not perform this step. However, the **back substitution** is disadvantageous compared to the **simple scaling** and the multiplication by the unitary \mathbf{V} in the SVD approach.

6.6. Rank-Deficient Least Square

The SVD of rank-deficient \mathbf{A} is given by

$$\mathbf{A} = \mathbf{U}_{\text{red}} \Sigma_{\text{red}} \mathbf{V}_{\text{red}}^H$$

with $\mathbf{U}_{\text{red}} \in \mathbb{C}^{m \times r}$, $\Sigma_{\text{red}} \in \mathbb{C}^{r \times r}$ which is diagonal with positive diagonal elements, sub unitary $\mathbf{V}_{\text{red}} \in \mathbb{C}^{n \times r}$ and $r = \text{rank}(\mathbf{A})$.

$$\Sigma_{\text{red}}^2 \mathbf{V}_{\text{red}}^H \mathbf{x} = \Sigma_{\text{red}} \mathbf{U}_{\text{red}}^H$$

The equation is underdetermined with more unknowns n than equations r . There are infinitely many solutions for \mathbf{x} exist. We rewrite the equation with

$$\mathbf{x}_{\text{LS}} = \arg\min_{\mathbf{x}} \|\mathbf{x}\|_2^2 \quad \text{s.t. } \Sigma_{\text{red}} \mathbf{V}_{\text{red}}^H \mathbf{x} = \mathbf{U}_{\text{red}}^H \mathbf{b}$$

For rank-deficient LS we obtain the solution:

$$\mathbf{x}_{\text{LS}} = \mathbf{V}_{\text{red}} \Sigma_{\text{red}}^{-1} \mathbf{U}_{\text{red}}^H \mathbf{b}$$

The residual is given by

$$\mathbf{r} = (\mathbf{I} - \mathbf{U}_{\text{red}} \mathbf{U}_{\text{red}}^H) \mathbf{b}$$

that is, the residual of the rank-deficient LS problem is equal to the part of the right-hand side \mathbf{b} lying in the orthogonal complement of $\text{range}(\mathbf{U}_{\text{red}}) = \text{range}(\mathbf{A})$.

6.7. Moore-Penrose Pseudoinverse

The generalization for all LS problems is given by the Moore-Penrose pseudoinverse with $\Sigma^\dagger = [\Sigma_1^\dagger, \mathbf{0}]$ we can rewrite \mathbf{A}^\dagger as

$$\mathbf{A}^\dagger = \mathbf{V} [\text{diag}(\frac{1}{\sigma_1}, \dots, \frac{1}{\sigma_r}, 0, \dots, 0), \mathbf{0}] \mathbf{U}^H$$

When \mathbf{A} is full-rank and square $\mathbf{A}^\dagger = \mathbf{A}^{-1}$

For full-rank and tall $\mathbf{A} \in \mathbb{C}^{m \times n}$, $m > n$: $\mathbf{A}^\dagger = (\mathbf{A}^H \mathbf{A})^{-1} \mathbf{A}^H$

For full-rank and wide matrix $\mathbf{A} \in \mathbb{C}^{m \times n}$, $m < n$: $\mathbf{A}^\dagger = \mathbf{A}^H (\mathbf{A} \mathbf{A}^H)^{-1}$

6.8. Rank-Revealing QR Decomposition

The application of the rank-revealing QR decomposition of $\mathbf{A} \in \mathbb{C}^{m \times n}$ with $m \geq n$ gives

$$\mathbf{A} \Pi = \mathbf{Q}_{\text{red}} \mathbf{R}_{\text{red}}$$

where $\mathbf{Q}_{\text{red}} \in \mathbb{C}^{m \times n}$ is sub unitary and

$$\mathbf{R}_{\text{red}} = \begin{bmatrix} \mathbf{R}_{11} & \mathbf{R}_{12} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \in \mathbb{C}^{n \times n}$$

The partition $\mathbf{R}_{11} \in \mathbb{C}^{r \times r}$ is upper-triangular and full-rank with $r = \text{rank}(\mathbf{A})$

In every step the **column** whose lower part, that has not yet been brought to triangular form and **exhibits maximum norm** in the current step is **permuted** to the first column of the part that has not yet been transformed to triangular form.

7. Condition of a Problem

7.1. Norm-Wise Condition Number

In abstract form, when solving a problem e.g., solving an equation system, the following map must be applied

$$\mathbf{f} : \mathbb{C}^m \rightarrow \mathbb{C}^n; \mathbf{x} \mapsto \mathbf{y} = \mathbf{f}(\mathbf{x})$$

A prominent example for such a problem is solving an equation system of the form

$$\mathbf{A} \mathbf{v} = \mathbf{b}$$

Interpreting the matrix \mathbf{A} constituting the equation system as parameter of the problem, the problem to be solved is:

$$\mathbf{f} : \mathbb{C}^m \rightarrow \mathbb{C}^n; \mathbf{b} \mapsto \mathbf{v} = \mathbf{f}(\mathbf{b})$$

Note that a different \mathbf{A}_2 leads to a different problem \mathbf{f}_2 .

The norm-wise condition number is given by:

$$\kappa_N(\mathbf{f}, \mathbf{x}) = \frac{\|\mathbf{J}(\mathbf{x})\| \cdot \|\mathbf{x}\|}{\|\mathbf{f}(\mathbf{x})\|}$$

where

$$\mathbf{J}(\mathbf{x}) = \frac{\partial \mathbf{f}(\mathbf{x})}{\partial \mathbf{x}^T}$$

The perturbed outcome of the problem \mathbf{f} can be bounded as follows

$$\frac{\|\mathbf{f}(\mathbf{x} + \Delta \mathbf{x}) - \mathbf{f}(\mathbf{x})\|}{\|\mathbf{f}(\mathbf{x})\|} \leq \kappa_N(\mathbf{f}, \mathbf{x}) \frac{\|\Delta \mathbf{x}\|}{\|\mathbf{x}\|} + o(\|\Delta \mathbf{x}\|)$$

- **well-conditioned:** If $\kappa_N(\mathbf{f}, \mathbf{x}) \approx 1$, (e.g., 1,100). For a well-conditioned problem, it is possible to find algorithms with small errors.
- **ill-conditioned:** If $\kappa_N(\mathbf{f}, \mathbf{x}) \gg 1$, (e.g., 10^6 , 10^{10}). When the problem is ill-conditioned, large errors must be expected from all algorithms.
- **ill-posed:** If $\kappa_N(\mathbf{f}, \mathbf{x}) = \infty$. In this case, the solution $\mathbf{f}(\mathbf{x})$ either does not exist or is not unique. It can be expected that no algorithm is able to give a reasonable result.

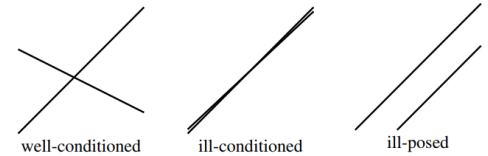


Fig. 6.3. Condition of intersecting two straight lines

$f(x)$	$J(x)$	κ_N
$f(x) = \sqrt{x}$	$J(x) = \frac{d\sqrt{x}}{dx} = \frac{1}{2} \frac{1}{\sqrt{x}}$	$\kappa_N = \frac{1}{2}$
$f(x) = ax$	$J(x) = a$	$\kappa_N = \frac{ a x }{ ax } = 1$
$f\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) = x_1 \cdot x_2$	$J(\mathbf{x}) = [x_2, x_1]$	$\kappa_N = \left \frac{x_1}{x_2} \right + \left \frac{x_2}{x_1} \right $

7.2. Component-wise Condition Number

The component-wise condition number is defined by

$$\begin{aligned} \kappa_C(f, \mathbf{x}) &= \frac{\|\mathbf{J}(\mathbf{x}) \odot \mathbf{x}^T\|_\infty}{\|\mathbf{f}(\mathbf{x})\|} \\ &= \frac{|\mathbf{J}^T(\mathbf{x})| |\mathbf{x}|}{\|\mathbf{f}(\mathbf{x})\|} \end{aligned}$$

the product $|\mathbf{J}^T(\mathbf{x})| |\mathbf{x}|$ is performed element wise.

The component-wise condition number is smaller than the norm-wise condition number if the infinity norm is employed. For $f : \mathbb{R}^m \rightarrow \mathbb{R}$ and when applying the infinity norm, the norm wise condition number is given by

$$\kappa_N(\mathbf{f}, \mathbf{x}) = \frac{\|\mathbf{J}(\mathbf{x})\|_\infty \cdot \|\mathbf{x}\|_\infty}{\|\mathbf{f}(\mathbf{x})\|}$$

as $f(\mathbf{x})$ is assumed to be scalar. Because $\mathbf{J}^T(\mathbf{x}) \in \mathbb{R}^{1 \times m}$ is a row vector for scalar $f(\mathbf{x})$, $\|\mathbf{J}^T(\mathbf{x})\|_\infty = \|\mathbf{J}(\mathbf{x})\|_1 = \sum_{i=1}^m |[\mathbf{J}(\mathbf{x})]_i|$ holds, and due to $\|\mathbf{x}\|_\infty = \max_{i \in \{1, \dots, m\}} |x_i|$

$$\kappa_N(f, \mathbf{x}) \geq \frac{|\mathbf{J}^T(\mathbf{x})| \cdot |\mathbf{x}|}{\|\mathbf{f}(\mathbf{x})\|} = \kappa_C(f, \mathbf{x})$$

f	$x_1 + x_2$	$x_1 - x_2$	$x_1 x_2$	$\frac{x_1}{x_2}$	$\frac{1}{x_1}$	$\sqrt{x_1}$
$\kappa_C(f; \mathbf{x})$	$\frac{ x_1 + x_2 }{ x_1+x_2 }$	$\frac{ x_1 - x_2 }{ x_1-x_2 }$	2	2	1	$\frac{1}{2}$

Table 6.1. Element-wise condition number of basic arithmetic operations

7.3. Condition of a Matrix

We can express the condition number of matrix $\mathbf{A} \in \mathbb{C}^{m \times n}$ with the help of the pseudoinverse \mathbf{A}^\dagger as

$$\kappa(\mathbf{A}) = \|\mathbf{A}\| \cdot \|\mathbf{A}^\dagger\|$$

When we employ the 2-norm we get the result

$$\kappa_2(\mathbf{A}) = \|\mathbf{A}\| \cdot \|\mathbf{A}^\dagger\| = \frac{\sigma_1}{\sigma_n}$$

Wide Matrix For a wide matrix we divide the vector \mathbf{x} into a component $\mathbf{x}_\perp^\mathbf{A}$ in the nullspace of \mathbf{A} and a component $\mathbf{x}_\mathbf{A}$ perpendicular to the nullspace of \mathbf{A} by employing the project $\mathbf{P}_\mathbf{A}$ onto the range of $\mathbf{A}^\mathbf{H}$, that is, $\mathbf{x}_\perp^\mathbf{A} = \mathbf{P}_\perp^\mathbf{A} \mathbf{x}$ and $\mathbf{x}_\mathbf{A} = \mathbf{P}_\mathbf{A} \mathbf{x}$, respectively, such that $\mathbf{x}_\mathbf{A}^\mathbf{H} \mathbf{x}_\perp^\mathbf{A} = 0$. Substituting into κ_N results in

$$\kappa_N(\mathbf{A} \mathbf{x}, \mathbf{x}) = \frac{\|\mathbf{A}\| \cdot \|\mathbf{x}_\mathbf{A} + \mathbf{x}_\perp^\mathbf{A}\|}{\|\mathbf{A} \mathbf{x}_\mathbf{A}\|}$$

7.3.1. Condition of Unitary Matrix

$$\kappa(\mathbf{U}) = \|\mathbf{U}\| \cdot \|\mathbf{U}^\mathbf{H}\| = 1$$

8. Stability of an Algorithm

The algorithm $\tilde{\mathbf{f}}(\mathbf{x})$ to solve a problem $\mathbf{f}(\mathbf{x})$ is denoted as

$$\tilde{\mathbf{f}}: \mathbb{C}^m \rightarrow \mathbb{C}^n; \mathbf{x} \rightarrow \tilde{\mathbf{y}} = \tilde{\mathbf{f}}(\mathbf{x})$$

8.1. Accuracy

For any possible value for the input \mathbf{x} , the relative error of the outcome $\tilde{\mathbf{f}}(\mathbf{x})$ of the algorithm is in the order of ϵ_M , which is the least error to be expected.

$$\forall \mathbf{x} \in \mathbb{C}^m: \frac{\|\tilde{\mathbf{f}}(\mathbf{x}) - \mathbf{f}(\mathbf{x})\|}{\|\mathbf{f}(\mathbf{x})\|} = \mathcal{O}(\epsilon_M)$$

8.2. Stability

An algorithm $\tilde{\mathbf{f}}$ is stable, if

$$\forall \mathbf{x} \in \mathbb{C}^m: \exists \tilde{\mathbf{x}} \text{ with } \frac{\|\tilde{\mathbf{x}} - \mathbf{x}\|}{\|\mathbf{x}\|} = \mathcal{O}(\epsilon_M)$$

such that

$$\frac{\|\tilde{\mathbf{f}}(\mathbf{x}) - \mathbf{f}(\tilde{\mathbf{x}})\|}{\|\mathbf{f}(\tilde{\mathbf{x}})\|} = \mathcal{O}(\epsilon_M)$$

It is not possible to infer the accuracy of an algorithm from the observation that it is stable.

8.3. Backward Stability

An algorithm $\tilde{\mathbf{f}}$ is backward stable, if

$$\forall \mathbf{x} \in \mathbb{C}^m: \exists \tilde{\mathbf{x}} \text{ with } \frac{\|\tilde{\mathbf{x}} - \mathbf{x}\|}{\|\mathbf{x}\|} = \mathcal{O}(\epsilon_M)$$

such that

$$\tilde{\mathbf{f}}(\mathbf{x}) = \mathbf{f}(\tilde{\mathbf{x}})$$

This definition of backward stability is stronger than the definition of stability before, since the relative error between the output of the algorithm for the true input \mathbf{x} , i.e., $\tilde{\mathbf{f}}(\mathbf{x})$, and the output of the underlying problem $\mathbf{f}(\tilde{\mathbf{x}})$ for the perturbed input $\tilde{\mathbf{x}}$ must be zero. **Backward stability implies stability.**

$$\frac{\|\tilde{\mathbf{f}}(\mathbf{x}) - \mathbf{f}(\mathbf{x})\|}{\|\mathbf{f}(\mathbf{x})\|} = \mathcal{O}(\kappa_N(\mathbf{f}, \mathbf{x}) \epsilon_M)$$

9. System of Equations

$$\mathbf{A} = \mathbf{L} \mathbf{U}$$

with the unit lower-triangular $\mathbf{L} \in \mathbb{C}^{m \times m}$ and the upper-triangular $\mathbf{U} \in \mathbb{C}^{m \times m}$. The coefficients for \mathbf{L} are given as $l_{i,j} = \frac{a_{ij}}{a_{jj}}$

$$\mathbf{L} = \mathbf{I} + l_1 \mathbf{e}_1^\mathbf{T} + \dots + l_{m-1} \mathbf{e}_{m-1}^\mathbf{T}$$

$$= \begin{bmatrix} 1 & & & \\ l_{21} & 1 & & \\ \vdots & \ddots & \ddots & \\ l_{m,1} & \dots & l_{m,m-1} & 1 \end{bmatrix}$$

Algorithm 13 Gaussian Elimination without Pivoting

```

1: function GAUSSIANELIMINATIONNOPIVOTING(A)      ▷ square full-rank  $\mathbf{A} \in \mathbb{C}^{m \times m}$ 
2:    $\mathbf{U} \leftarrow \mathbf{A}$ 
3:    $\mathbf{L} \leftarrow \mathbf{I}$                                      ▷ initialize diagonal of  $\mathbf{L}$ 
4:   for  $k = 1$  to  $m - 1$  do
5:      $[L]_{k+1:m,k} \leftarrow \frac{1}{u_{kk}} [U]_{k+1:m,k}$       ▷ coefficients for Gaussian elimination step  $k$ 
6:      $[U]_{k+1:m,k} \leftarrow 0$                          ▷ elimination for column  $k$ 
7:      $[U]_{k+1:m,k+1:m} \leftarrow [U]_{k+1:m,k+1:m} - [L]_{k+1:m,k} [U]_{k,k+1:m}$   ▷ elimination for rest
8:   end for
9:   return  $\mathbf{L}, \mathbf{U}$ 
10: end function

```

Operation Count Algorithm 13: #FLOPS = $\mathcal{O}(\frac{2m^3}{3})$

Algorithm 14 Solving Equation System with Gaussian Elimination

```

1: function SOLVINGEQUATIONSYSTEMLU(A, b)
2:    $\mathbf{L}, \mathbf{U} \leftarrow \text{lu}(\mathbf{A})$                          ▷ LU factorization of  $\mathbf{A}$ 
3:    $\mathbf{w} \leftarrow \mathbf{L}^{-1} \mathbf{b}$                              ▷ back substitution
4:    $\mathbf{x} \leftarrow \mathbf{U}^{-1} \mathbf{w}$                              ▷ back substitution
5:   return  $\mathbf{x}$ 
6: end function

```

Operation Count Algorithm 14: #FLOPS = $\mathcal{O}(\frac{2m^3}{3} + 2m^2)$

Worst case error for Gauss with LU is:

$\mathcal{O}(\rho_{\text{worst-case}} \cdot \epsilon_M)$ with $\rho_{\text{worst-case}} = 2^{m-1}$. Solving an equation system with QR does not have such problems.

9.1. Stability of Gaussian Elimination

Without pivoting it holds that

$$\tilde{\mathbf{L}} \tilde{\mathbf{U}} = \mathbf{A} + \Delta \mathbf{A} \quad \text{with} \quad \frac{\|\Delta \mathbf{A}\|}{\|\mathbf{L}\| \|\mathbf{U}\|} = \mathcal{O}(\epsilon_M)$$

Gaussian elimination without pivoting is not backward stable!

9.2. Complete Pivoting

Leads to smallest coefficients for back substitution procedure. We have to compare $(m - k + 1)^2$ elements in the k -th step. All entries in \mathbf{L} below the diagonal are smaller than one.

Operation Count: #FLOPS = $\mathcal{O}(\frac{m^3}{3})$

Full pivoting ensures that all off-diagonal elements of the resulting \mathbf{U} are smaller than or equal to the diagonal entry in the same row. Thus complete pivoting also improves the properties of the second back substitution.

Algorithm 15 Gaussian Elimination with Partial Pivoting

```

1: function GAUSSIANELIMINATIONPIVOTING(A)
2:    $\mathbf{U} \leftarrow \mathbf{A}, \mathbf{L} \leftarrow \mathbf{I}, \mathbf{I} \leftarrow 1$            ▷ initialization
3:   for  $k = 1$  to  $m - 1$  do
4:      $i_{\max} \leftarrow \text{MAX}_{k \leq i \leq m} |u_{ik}|$            ▷ find maximum pivot
5:      $[U]_{k,k:m} \leftrightarrow [U]_{i_{\max},k:m}$                  ▷ exchange rows  $i_{\max}$  and  $k$  of  $\mathbf{U}$ 
6:      $[L]_{k,1:k-1} \leftrightarrow [L]_{i_{\max},1:k-1}$              ▷ exchange rows of lower non-zero part of  $\mathbf{L}$ 
7:      $[\mathbf{I}]_k \leftrightarrow [\mathbf{I}]_{i_{\max}}$                          ▷ include permutation in  $\mathbf{I}$ 
8:      $[L]_{k+1:m,k} \leftarrow \frac{1}{u_{kk}} [U]_{k+1:m,k}$        ▷ coefficients for Gaussian elimination step  $k$ 
9:      $[U]_{k+1:m,k} \leftarrow 0$                          ▷ elimination for column  $k$ 
10:     $[U]_{k+1:m,k+1:m} \leftarrow [U]_{k+1:m,k+1:m} - [L]_{k+1:m,k} [U]_{k,k+1:m}$   ▷ elimination of rest
11:  end for
12:  return  $\mathbf{L}, \mathbf{U}, \mathbf{I}$ 
13: end function

```

Operation Count Algorithm 15: #FLOPS = $\mathcal{O}(\frac{2m^3}{3} + \frac{m^2}{2})$

9.3. Partial Pivoting

With partial pivoting we have

$$\mathbf{\Pi A} = \mathbf{L} \mathbf{U}$$

For partial pivoting $m - k + 1$ elements must be compared in step k

Operation Count: #FLOPS = $\mathcal{O}(\frac{m^2}{2})$

The computational costs of partial pivoting compared to the LU factorization ($\mathcal{O}(\frac{2m^3}{3})$) are neglectable.

Gaussian elimination with partial pivoting is potentially backward stable. With growth factor $\rho = \frac{\max_{i,j} (|u_{i,j}|)}{\max_{i,j} |a_{i,j}|}$

With pivoting it holds: $\tilde{\mathbf{L}} \tilde{\mathbf{U}} = \mathbf{A} + \Delta \mathbf{A}$ with $\frac{\|\Delta \mathbf{A}\|}{\|\mathbf{L}\| \|\mathbf{U}\|} = \mathcal{O}(\rho \epsilon_M)$

9.4. Cholesky Factorization

A positive definite matrix \mathbf{A} is Hermitian i.e. $\mathbf{A}^\mathbf{H} = \mathbf{A}$. All eigenvalues of a Hermitian \mathbf{A} are real-valued. If \mathbf{A} is not only Hermitian but positive definite, all eigenvalues are not only real-valued but also positive!

$$\forall \mathbf{x} \in \mathbb{C}^m: \mathbf{x}^\mathbf{H} \mathbf{A} \mathbf{x} > 0$$

$$\mathbf{A} \succ 0$$

For Cholesky the matrix \mathbf{A} must be positive definite!

$$\mathbf{A} \succ 0$$

All diagonal elements of a positive definite matrix are greater than zero!

$$\mathbf{e}_k^\mathbf{T} \mathbf{A} \mathbf{e}_k^\mathbf{T} = [\mathbf{A}]_{k,k} > 0$$

With the Cholesky factorization the root of matrix can be defined as

$$\mathbf{A}^{1/2} = \mathbf{L} \mathbf{D}^{1/2}$$

$$\mathbf{A} = \mathbf{A}^{1/2} \mathbf{A}^{\mathbf{H}/2}$$

Algorithm 16 Cholesky Factorization

```

1: function CHOLESKYFACTORIZATION(A)
2:    $\mathbf{L} \leftarrow \mathbf{A}$ 
3:   for  $k = 1$  to  $m - 1$  do
4:     for  $j = k + 1$  to  $m$  do
5:        $L_{j:m,j} \leftarrow L_{j:m,j} - \frac{L_{jk}}{L_{kk}} L_{j:m,k}$    ▷ elimination only of lower triangular part
6:     end for
7:   end for
8:    $\mathbf{D} \leftarrow \text{diag}(\text{diag}(\mathbf{L}))$ 
9:    $\mathbf{L} \leftarrow \text{tril}(\mathbf{L} \mathbf{D}^{-1})$                        ▷ copy the lower triangular part, set rest to zero
10:  return  $\mathbf{L}, \mathbf{D}$ 
11: end function

```

Operation Count Algorithm 16: #FLOPS = $\mathcal{O}(\frac{1}{3} m^3)$

$$\begin{aligned} \# \text{FLOP}_k &= \sum_{j=k+1}^m 2m - 2j + 3 = \sum_{j'=1}^{m-k} 2j' + 1 \\ &= (m - k)(m - k + 1) + m - k \approx m^2 - 2mk + k^2 \end{aligned}$$

$$\begin{aligned} \# \text{FLOPS} &= \sum_{k=1}^{m-1} \# \text{FLOPS}_k = \sum_{k=1}^{m-1} m^2 - 2mk + k^2 \\ &= (m - 1)m^2 - m^2(m - 1) + \frac{(2m - 1)m(m - 1)}{6} \end{aligned}$$

9.5. Stability of Cholesky Factorization

Contrary to Gaussian Elimination, the Cholesky factorization can be used without pivoting. If, however, the stability of the Cholesky factorization must be improved, pivoting can be employed.

$$\mathbf{\Pi A} \mathbf{\Pi}^\mathbf{T} = \mathbf{L} \mathbf{D} \mathbf{L}^\mathbf{H}$$

Cholesky factorization, even without pivoting, can be shown to be backward stable

9.6. Hessenberg-Form

A matrix whose elements below the first lower sub-diagonal are all zero is called **Hessenberg matrix** \mathbf{H}

$$\mathbf{Q}^\mathbf{H} \mathbf{A} \mathbf{Q} = \begin{bmatrix} \times & \times & \dots & \times & \times \\ \times & \times & \dots & \times & \times \\ 0 & \times & \dots & \times & \times \\ \vdots & \ddots & \ddots & \vdots & \vdots \\ 0 & \dots & 0 & \times & \times \end{bmatrix}$$

The Hessenberg-Form can be achieved for a matrix $\mathbf{A} \in \mathbb{C}^{m \times m}$ if we employ unitary reflection matrices \mathbf{Q}_i with the resulting reflection matrix $\mathbf{Q} = \mathbf{Q}_1 \cdots \mathbf{Q}_{m-2}$.

$$\mathbf{Q} = \begin{bmatrix} \mathbf{I}_i & \mathbf{0}^\mathbf{T} \\ \mathbf{0} & \mathbf{F}_i \end{bmatrix} \in \mathbb{C}^{m \times m}$$

With the Householder reflectors $\mathbf{F}_i \in \mathbb{C}^{m-i+1 \times m-i+1}$ and $\mathbf{I}_i \in \mathbb{C}^{i-1 \times i-1}$

For Hermitian \mathbf{A} , i.e. $\mathbf{A}^\mathbf{H} = \mathbf{A}$ the result of the transformations leading to a Hessenberg matrix is a Hermitian Hessenberg matrix, that is, a tri-diagonal matrix.

$$\mathbf{H} = \mathbf{Q}^\mathbf{H} \mathbf{A} \mathbf{Q} = \begin{bmatrix} \times & \times & 0 & 0 & 0 & \dots & 0 \\ \times & \times & \times & 0 & 0 & \dots & 0 \\ 0 & \times & \times & \times & 0 & \dots & 0 \\ \vdots & \ddots & \ddots & \ddots & \ddots & \ddots & \vdots \\ 0 & \dots & 0 & \times & \times & \times & 0 \\ 0 & \dots & 0 & 0 & \times & \times & \times \\ 0 & \dots & 0 & 0 & 0 & \times & \times \end{bmatrix} = \mathbf{H}^\mathbf{H}$$