

Numerical Methods for Computational Science and Engineering

Fall Semester 2017 (HS17)

Prof. Rima Alaifari, SAM, ETH Zurich

Note on stability: $\kappa_{A^T A} = \kappa_A^2$ $\kappa_A = \text{cond}(A)$
condition number squares

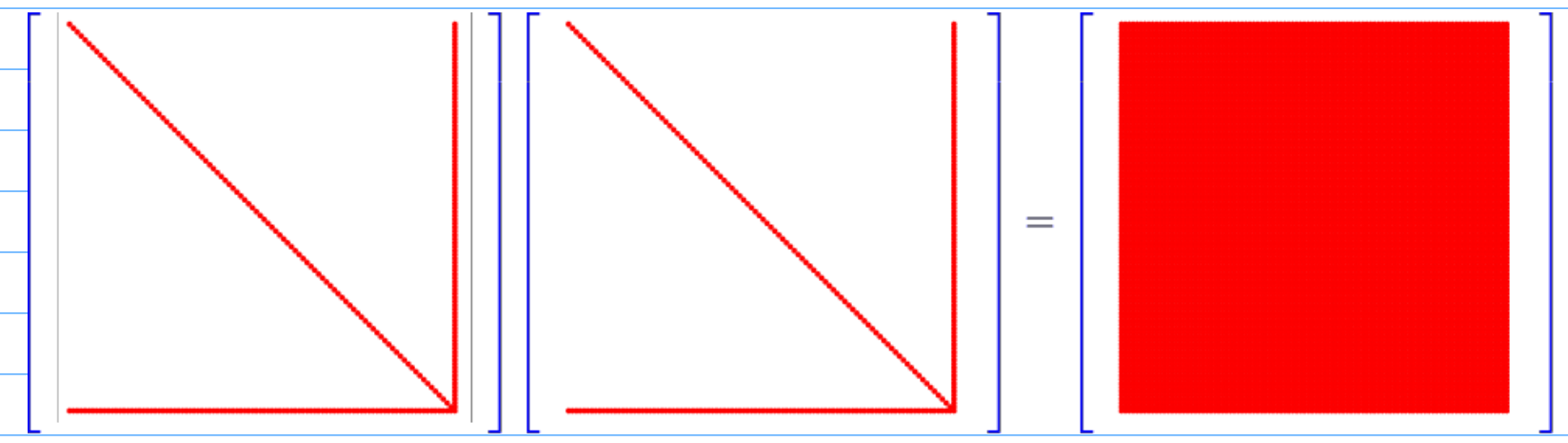
Example (3.2.4): $A = \begin{bmatrix} 1 & 1 \\ \delta & 0 \\ 0 & \delta \end{bmatrix}$ $A^T A = \begin{bmatrix} 1+\delta^2 & 1 \\ 1 & 1+\delta^2 \end{bmatrix}$

If $\delta \approx \sqrt{\text{EPS}}$ e.g. $\delta = \frac{\sqrt{\text{EPS}}}{2}$: $1+\delta^2 \approx 1 + \frac{\text{EPS}}{4} \underset{\uparrow}{=} 1$
in machine number arithmetic

As an element of $M^{2,2}$, $A^T A$ is not regular.

Further note:

$A \text{ sparse} \not\Rightarrow A^T A \text{ sparse}$



arrow matrix A: sparse \checkmark $A^T A$ dense

- ① Squaring condition number
 - ② Loss of sparsity
- } challenges

Extended normal equation: (3.2.7) ← maintain sparsity

$$r := Ax - b$$

$$A^H Ax = A^H b \Leftrightarrow B \begin{bmatrix} r \\ x \end{bmatrix} := \underbrace{\begin{bmatrix} -I & A \\ A^H & 0 \end{bmatrix}}_B \begin{bmatrix} r \\ x \end{bmatrix} = \begin{bmatrix} b \\ 0 \end{bmatrix}$$

$$-r + Ax = b \quad (I)$$

$$A^H r = A^H (Ax - b) = 0 \quad (II)$$

if A is sparse $\Rightarrow B$ is also sparse

BUT: conditioning is not improved

More general form: $r := \alpha^{-1}(Ax - b)$

for some choice of parameter $\alpha > 0$

$$A^H Ax = A^H b \Leftrightarrow B_\alpha \begin{bmatrix} r \\ x \end{bmatrix} := \begin{bmatrix} -\alpha I & A \\ A^H & 0 \end{bmatrix} \begin{bmatrix} r \\ x \end{bmatrix} = \begin{bmatrix} b \\ 0 \end{bmatrix}$$

Good choice of α can give better C_{B_α}
hopefully $C_{B_\alpha} \approx C_A$
(instead of $C_B \approx C_{A^H A}$)

[Note: this is not regularization

$$(A^H A + \alpha I)x = A^H b$$

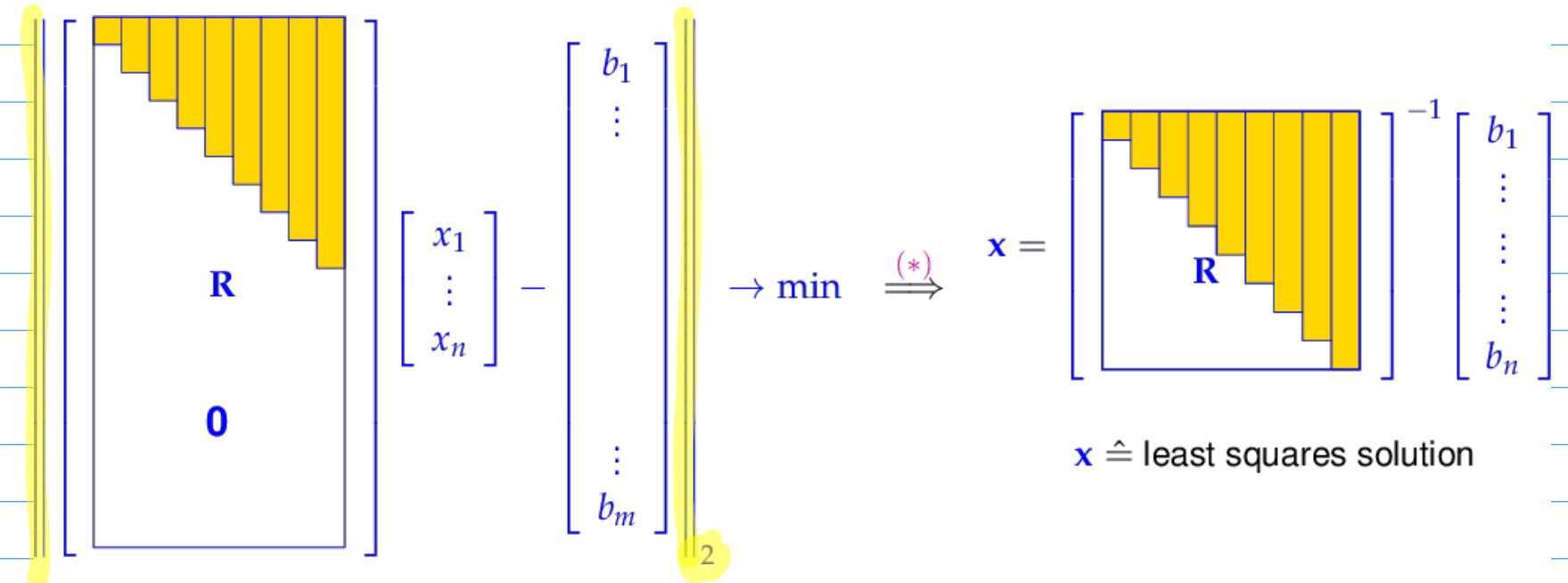
$$r = Ax - b \quad \tilde{B}_\alpha = \begin{bmatrix} -I & A \\ A^H & \alpha I \end{bmatrix}$$

3.3. Orthogonal Transformation Methods

Considers least-sq. problem $Ax=b$ $A \in \mathbb{R}^{m,n}$
 $m \gg n$

and A has full rank $\text{rank}(A)=n$

Idea: instead of solving $Ax=b$
find \tilde{A}, \tilde{b}
with $\text{lsq}(\tilde{A}, \tilde{b}) = \text{lsq}(A, b)$ s.t.
 $\tilde{A}x = \tilde{b}$ is easier to solve.
↑
triangular system matrix



R will be regular if $\text{rank}(A)=n$.

Idea: If we have a (transformation) matrix $T \in \mathbb{R}^{m,m}$ satisfying

$$\|Ty\|_2 = \|y\|_2 \quad \forall y \in \mathbb{R}^m, \tag{3.3.1}$$

then $\underset{y \in \mathbb{R}^n}{\text{argmin}} \|Ay - b\|_2 = \underset{y \in \mathbb{R}^n}{\text{argmin}} \|\tilde{A}y - \tilde{b}\|_2$,

where $\tilde{A} = TA$ and $\tilde{b} = Tb$.

\Rightarrow multiply our LSE with a matrix $T \in \mathbb{R}^{m,m}$
 with property $\|Ty\|_2 = \|y\|_2 \quad \forall y \in \mathbb{R}^m$
 yields modified LSE with set lsq
 unchanged!

\leadsto orthogonal / unitary matrices

Definition 3.3.4. Unitary and orthogonal matrices \rightarrow [?, Sect. 2.8]

- $Q \in \mathbb{K}^{n,n}$, $n \in \mathbb{N}$, is **unitary**, if $Q^{-1} = Q^H$.
- $Q \in \mathbb{R}^{n,n}$, $n \in \mathbb{N}$, is **orthogonal**, if $Q^{-1} = Q^T$.

$$Q^H Q = I = Q Q^H$$

Theorem 3.3.5. Preservation of Euclidean norm

A matrix is unitary/orthogonal, if and only if the associated linear mapping preserves the 2-norm:

$$Q \in \mathbb{K}^{n,n} \text{ unitary} \Leftrightarrow \|Qx\|_2 = \|x\|_2 \quad \forall x \in \mathbb{K}^n.$$

Goal: Transform $Ax = b$ (*) to

$$Q^T A x = Q^T b \quad (**)$$

\uparrow orthogonal

s.t. $Q^T A$ is upper triangular

Solving $A^T A x = A^T b$ vs. solving

$Q^T A x = Q^T b$ in least-sq. sense

$$A^T \underbrace{Q Q^T}_{=I} A x = A^T \underbrace{Q Q^T}_{=I} b$$

$$\Leftrightarrow \underline{A^T A x = A^T b}$$

Problems (*) and (**) are equivalent in the least-squares sense

Recall: $A \in \mathbb{R}^{m,n}$, $Q \in \mathbb{R}^{m,m}$ orthogonal

if $R := Q^T A \in \mathbb{R}^{m,n}$ was upper triangular

least-squares problem becomes

$$\underbrace{Q^T A}_R x = Q^T b$$

$$(**) \Leftrightarrow Rx = Q^T b$$

as cheap as back-substitution

If $A \in \mathbb{R}^{m,n}$ can be decomposed in

$$A = QR$$

↑ ↑
orth. $m \times m$ triang. $m \times n$

$$A^T A x = A^T b \Leftrightarrow$$

$$R^T \underbrace{Q^T Q}_I R x = R^T Q^T b$$

$$\Leftrightarrow R^T R x = R^T Q^T b$$

$$\begin{aligned} \underline{x} &= (R^T R)^{-1} R^T Q^T b \\ &= \underline{R^{-1} Q^T b} \end{aligned}$$

Instead of solving $A^T A x = A^T b$ with

$$\kappa_{A^T A} = \kappa_A^2$$

solving for $Rx = Q^T b$ is better conditioned

$$\text{with } \kappa_R = \kappa_A$$

Next question: Does such a QR decomposition always exist? How to find it?

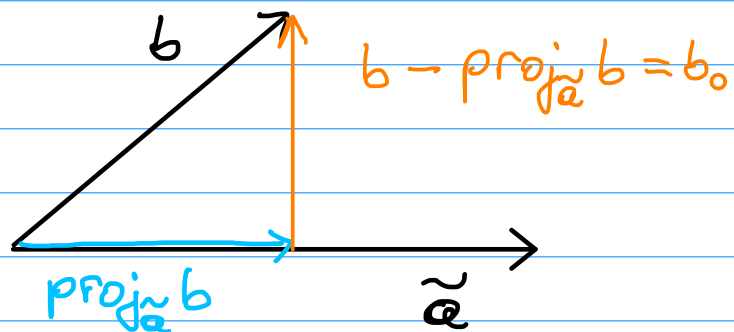
3.3.3 QR-Decomposition

First approach: Gram-Schmidt orthogonalization

Orthogonalization of 2 linearly independent

vectors $a, b \in \mathbb{R}^m$

$$\tilde{a} := \frac{a}{\|a\|_2}$$



$\text{proj}_{\tilde{a}} b$ e.g. as least-sq. problem:

find τ that minimizes $\|\tau \cdot \tilde{a} - b\|_2$

$$\Leftrightarrow \underbrace{\tilde{a}^T \tilde{a}}_{=1} \tau - \tilde{a}^T b = 0$$

$$\tau = \langle \tilde{a}, b \rangle$$

$$\text{proj}_{\tilde{a}} b = \langle \tilde{a}, b \rangle \tilde{a} \quad (= (\tilde{a}^T b) \tilde{a})$$

$$b_0 = b - \text{proj}_{\tilde{a}} b = b - \langle \tilde{a}, b \rangle \tilde{a} \quad (= b - \frac{\langle a, b \rangle}{\|a\|_2^2} a)$$

$$\tilde{b} := \frac{b_0}{\|b_0\|_2}$$

ONS: $\{\tilde{a}, \tilde{b}\}$

Gram-Schmidt: Procedure for set of k lin. ind. vectors

- 1: $q^1 := \frac{a^1}{\|a^1\|_2}$ % 1st output vector
- 2: for $j = 2, \dots, k$ do
 - { % Orthogonal projection
 - 3: $q^j := a^j$
 - 4: for $l = 1, 2, \dots, j-1$ do (GS)
 - 5: $\{ q^j \leftarrow q^j - \langle a^j, q^l \rangle q^l \}$
 - 6: if $(q^j = 0)$ then STOP
 - 7: else $\{ q^j \leftarrow \frac{q^j}{\|q^j\|_2} \}$
 - 8: }

Theorem 3.3.7. Span property of G.S. vectors

If $\{a^1, \dots, a^k\}$ is linearly independent, then Algorithm (GS) computes orthonormal vectors q^1, \dots, q^k satisfying

$$\text{Span}\{q^1, \dots, q^l\} = \text{Span}\{a^1, \dots, a^l\}, \quad (1.5.2)$$

for all $l \in \{1, \dots, k\}$.

Given $\{a^1, \dots, a^k\}$ lin. ind. \rightarrow output $\{q^1, \dots, q^k\}$
 orth. normal system
 (ONS)

For matrix $A = [a^1 \ a^2 \ \dots \ a^n] \in \mathbb{R}^{m,n}$

Step 1:

$$\begin{pmatrix} | & | & & | \\ a^1 & a^2 & \dots & a^n \\ | & | & & | \end{pmatrix} \begin{pmatrix} \tilde{t}_{11} & & & \\ & 1 & & \\ & & \ddots & \\ & & & 1 \end{pmatrix} = \begin{pmatrix} | & | & & | \\ \tilde{t}_{11} a^1 & a^2 & \dots & a^n \\ | & | & & | \end{pmatrix}$$

Step 2:

$$\begin{pmatrix} | & | & & | \\ \tilde{t}_{11} a^1 & a^2 & \dots & a^n \\ | & | & & | \end{pmatrix} \begin{pmatrix} 1 & \tilde{t}_{12} & 0 & \\ \theta & \tilde{t}_{22} & 1 & \\ \vdots & \theta & & \ddots \\ \theta & \theta & & \theta & \ddots & 1 \end{pmatrix} =$$

$$\begin{matrix} \tilde{t}_{12} = \tilde{t}_{11} \tilde{t}_{12} \\ \nearrow \\ = \end{matrix} \begin{pmatrix} | & | & & | \\ \tilde{t}_{11} a^1 & \tilde{t}_{12} a^1 + \tilde{t}_{22} a^2 & \dots & a^n \\ | & | & & | \end{pmatrix}$$

Altogether: Series of multiplications from the
 right by upper triangular matrices

$$Q = A \underbrace{T_1 T_2 \dots T_n}_{=: T}$$

$$Q = [q^1 \dots q^n] \in \mathbb{R}^{m,n} \quad \text{s.t.}$$

$$q^1 = t_{11} a^1$$

$$q^2 = t_{12} a^1 + t_{22} a^2$$

$$q^3 = t_{13} a^1 + t_{23} a^2 + t_{33} a^3$$

⋮

$$q^n = t_{1n} a^1 + \dots + t_{nn} a^n$$

and $Q = AT$

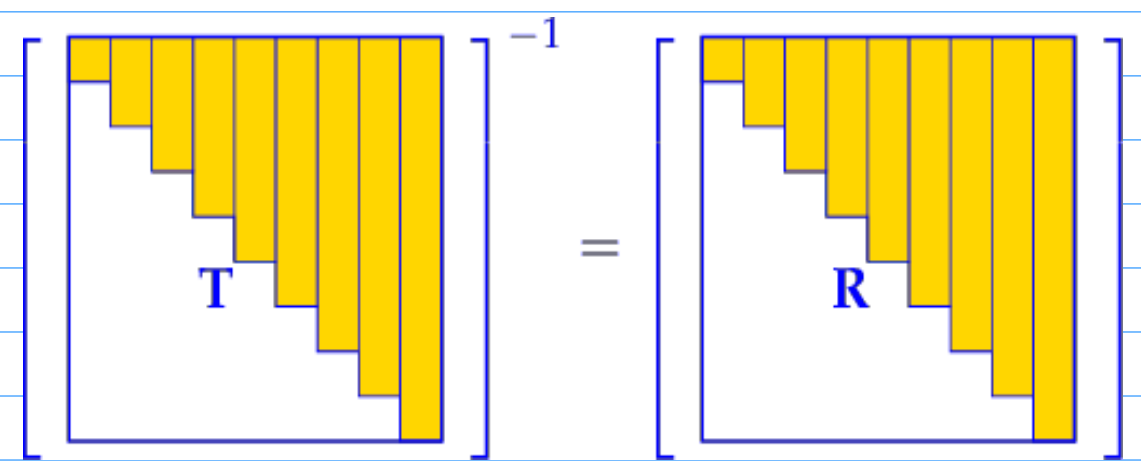
$$T \in \mathbb{R}^{n,n} \quad T = (t_{ij})_{i,j=1}^n$$

upper triangular

T is regular because $\{a^1, \dots, a^n\}$
and $\{q^1, \dots, q^n\}$ are lin. ind.

T reg. upper triangular

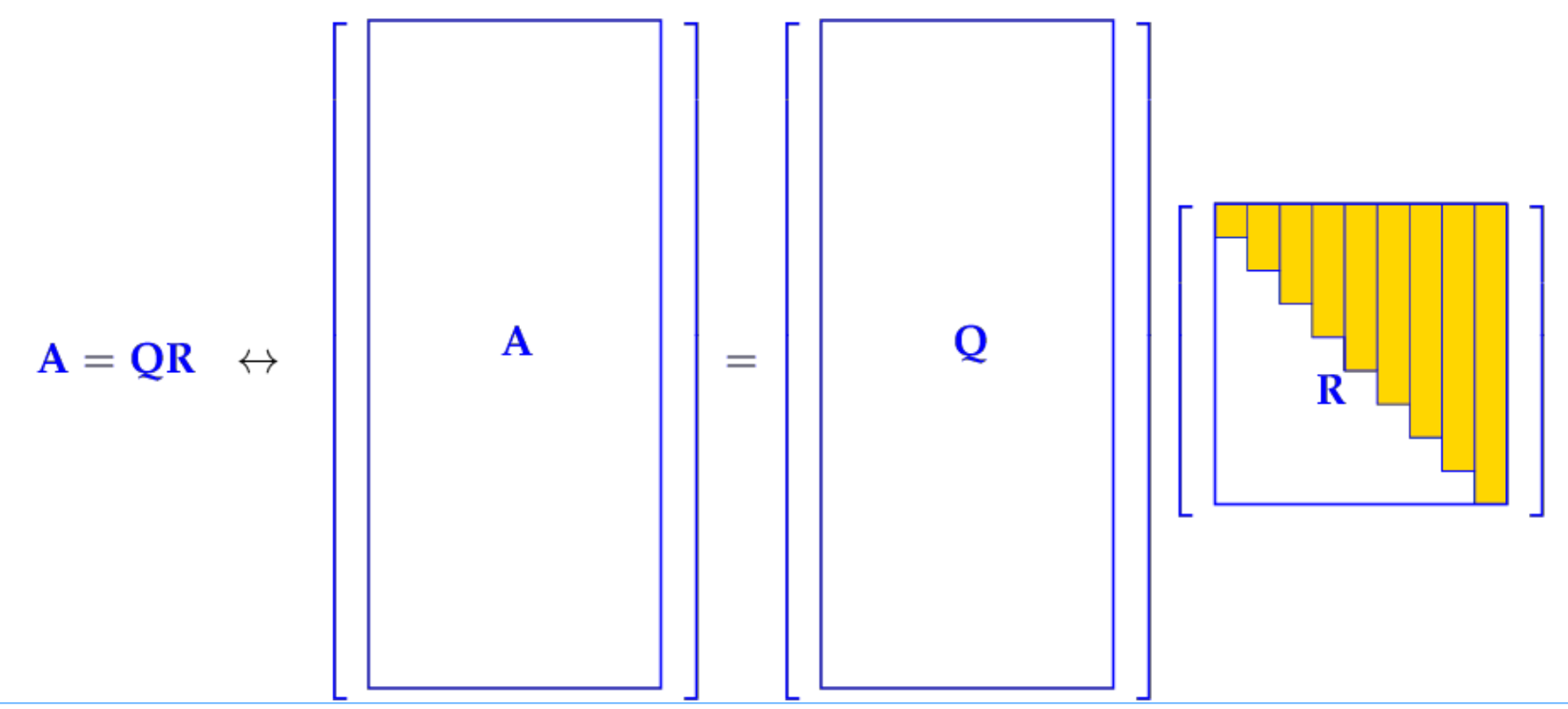
$\Rightarrow T^{-1}$ reg. upper triangular



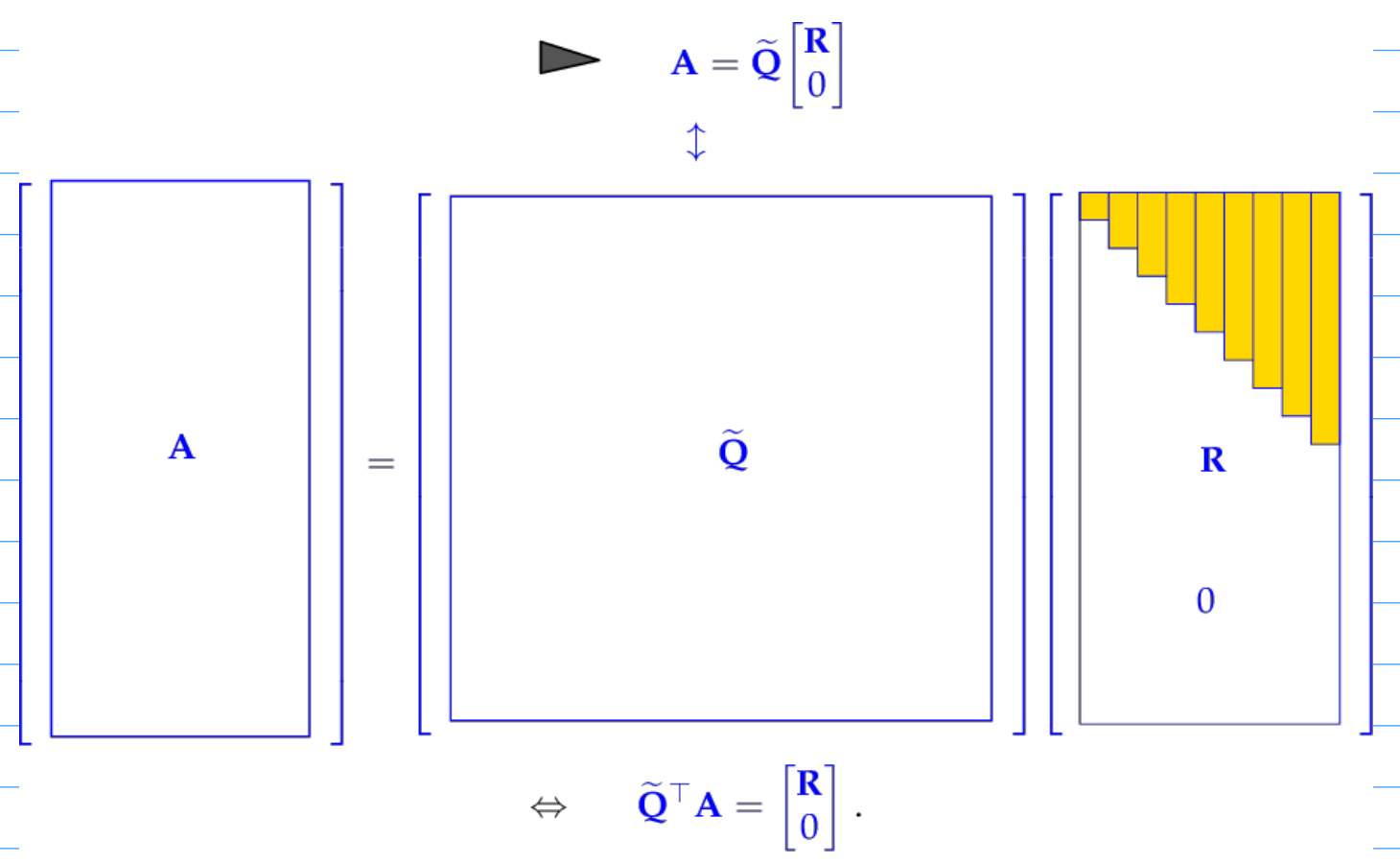
$$R := T^{-1}$$

$\Rightarrow A = QR$

Thus, by (3.3.8), we have found an *upper triangular* $R := T^{-1} \in \mathbb{R}^{n,n}$ such that



Argument:



[add $m-n$ zero rows to R
 and add $m-n$ columns to Q s.t. \tilde{Q}
 is an orth. (unitary) matrix]

Theorem 3.3.9. QR-decomposition → [?, Satz 5.2], [?, Sect. 7.3]

For any matrix $A \in \mathbb{K}^{n,k}$ with $\text{rank}(A) = k$ there exists

- (i) a unique Matrix $Q_0 \in \mathbb{R}^{n,k}$ that satisfies $Q_0^H Q_0 = I_k$, and a unique *upper triangular* Matrix $R_0 \in \mathbb{K}^{k,k}$ with $(R_0)_{i,i} > 0, i \in \{1, \dots, k\}$, such that

$$A = Q_0 \cdot R_0 \quad (\text{"economical" QR-decomposition}),$$

- (ii) a *unitary* Matrix $Q \in \mathbb{K}^{n,n}$ and a unique *upper triangular* $R \in \mathbb{K}^{n,k}$ with $(R)_{i,i} > 0, i \in \{1, \dots, n\}$, such that

$$A = Q \cdot R \quad (\text{full QR-decomposition}).$$

If $\mathbb{K} = \mathbb{R}$ all matrices will be real and Q is then *orthogonal*.

$$A = Q_0 R_0, \quad Q_0 \in \mathbb{K}^{n,k}, \quad R_0 \in \mathbb{K}^{k,k} \text{ upper triangular,}$$

$$A = Q_0 R_0 \quad (3.3.10)$$

$Q_0^H Q_0 = \hat{I}_k$ (Q₀ has orthonormal columns)

$$Q^H Q = I_n = Q Q^H$$

$$A = QR, \quad Q \in \mathbb{K}^{n,n}, \quad R \in \mathbb{K}^{n,k},$$

$$A = QR \quad (3.3.11)$$

Corollary 3.3.12. Uniqueness of QR-factorization

The "economical" QR-factorization (3.3.3.1) of $A \in \mathbb{K}^{m,n}, m \geq n$, with $\text{rank}(A) = n$ is unique, if we demand $(R_0)_{ii} > 0$.

Proof: Suppose $Q_1 R_1 = A = Q_2 R_2$

$$Q_1 R_1 = Q_2 R_2$$

$$\Rightarrow Q_1 = Q_2 R_2 R_1^{-1} =: R \quad (\text{again upper triang.})$$

$$\Rightarrow \underline{I_k} = Q_1^H Q_1 = R^H \underbrace{Q_2^H Q_2}_{= I_k} R = \underline{R^H R}$$

R is upper triang. & orth.

$$R = \begin{pmatrix} \textcircled{r_{11}} \neq 0 & r_{12} & r_{13} & \dots & r_{1k} = 0 \\ & r_{22} & r_{23} & \dots & r_{2k} = 0 \\ & & \ddots & & \vdots \\ & & & & r_{kk} \end{pmatrix}$$

orth.: $\langle r_1, r_j \rangle = 0 \quad j \neq 1$

$$\Rightarrow r_{12} = r_{13} = \dots = r_{1k} = 0$$

$$\Rightarrow r_{22} \neq 0; \text{ using } \langle r_2, r_j \rangle = 0 \quad j \neq 2$$

ind.

$\Rightarrow R$ is a diag.

$$\Rightarrow R = \begin{pmatrix} r_{11} & & & & \\ & \ddots & & & \\ & & \ddots & & \\ & & & \ddots & \\ & & & & r_{kk} \end{pmatrix}$$

and $|r_{ii}|^2 = 1 \quad (R^H R = I_k)$

if $r_{ii} > 0 \Rightarrow R = I_k \quad (R = R_2 R_1^{-1})$

$$\Rightarrow R_2 = R_1 \quad \text{and} \quad Q_1 = Q_2$$

Note: Gram-Schmidt orthogonalization suffers from numerical instabilities [possible cancellation in subtraction \leadsto then dividing by ≈ 0]

Example: $a_1 = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$ $a_2 = \begin{pmatrix} 1+\varepsilon \\ 1 \end{pmatrix}$ $\varepsilon \ll 1$

$$\text{span}\{a_1, a_2\} = \text{span}\left\{\begin{pmatrix} 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \end{pmatrix}\right\}$$

but Gram-Schmidt gives

$$q_1 = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 \\ 1 \end{pmatrix} \quad \text{proj}_{q_1} a_2 = \langle a_2, q_1 \rangle q_1$$

$$= \frac{2+\varepsilon}{2} \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

subtraction of close-by numbers

$$\tilde{q}_2 = \begin{pmatrix} 1+\varepsilon \\ 1 \end{pmatrix} - \frac{2+\varepsilon}{2} \begin{pmatrix} 1 \\ 1 \end{pmatrix} = \frac{\varepsilon}{2} \begin{pmatrix} 1 \\ -1 \end{pmatrix}$$

$$\|\tilde{q}_2\|_2 = \frac{\varepsilon}{\sqrt{2}}$$

$$q_2 = \frac{\tilde{q}_2}{\|\tilde{q}_2\|_2} \leftarrow \text{division by } \varepsilon!$$

Numerically stable QR decomposition?

3.3.3.2 Computation of QR decomposition

Idea: Before: manipulation of columns of A
by multiplication from the right
[with triangular matrices]

$$Q = A T_1 \dots T_n$$

Now: Instead: Find a series of orthogonal transformations s.t. applied from the left yield triangular matrix:

$$\underbrace{Q_n \dots Q_2 Q_1}_{\text{row operations}} A = R$$

Then $Q := Q_1^T \dots Q_n^T$ yields $A = QR$.

→ idea similar to Gauss elimination
but now with orthogonal traps

Householder Transformations

Set of orthogonal transformations is rich

Intuition: orth. traps preserve

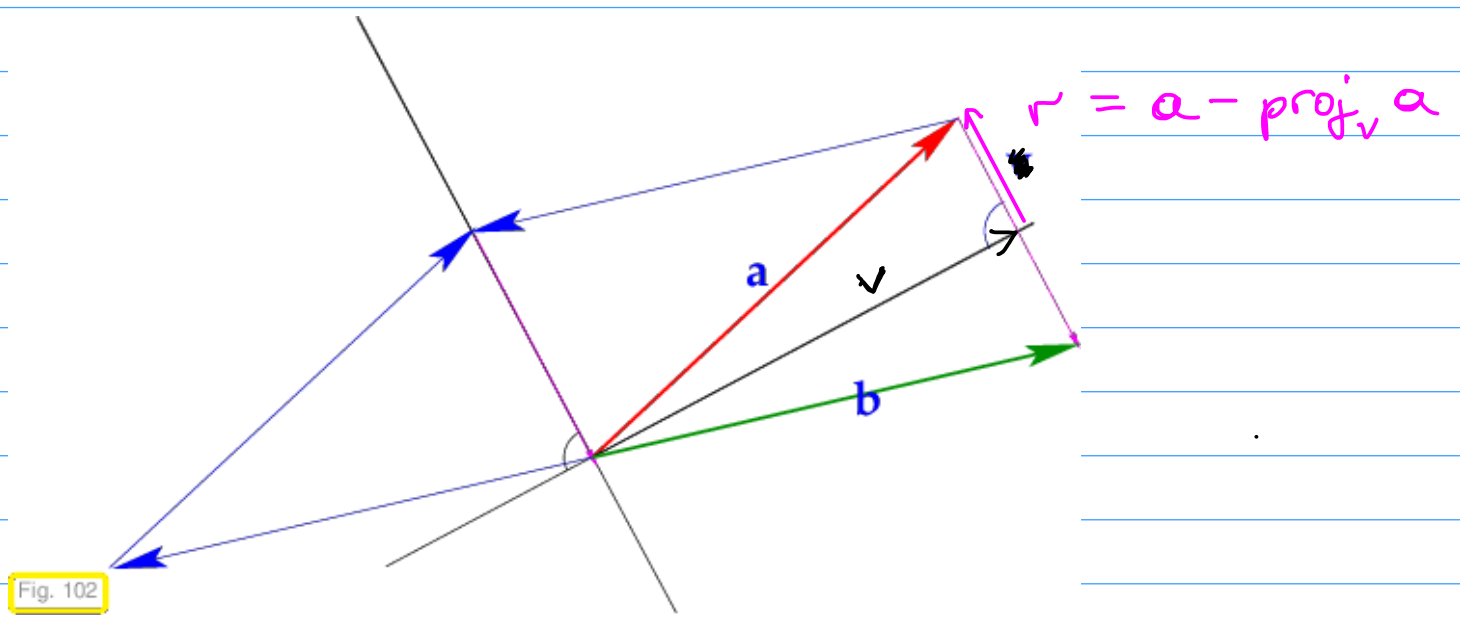
- lengths of vectors
- angles between vectors

→ can only rotate & reflect vectors

Idea: Use only reflections → can be represented by projections

Given a vector $a \in \mathbb{R}^m$ & reflected over a vector $v \in \mathbb{R}^m$

reflection: $b \in \mathbb{R}^m$



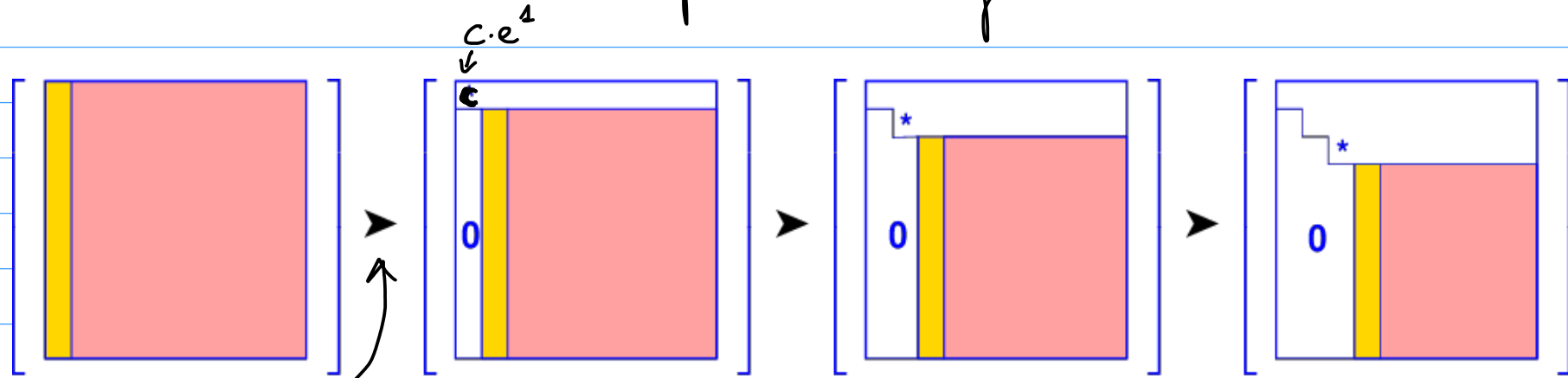
residual r is perpendicular to v

$$b = a - 2r = a - 2(a - \text{proj}_v a) = -a + 2 \text{proj}_v a = -(a - 2 \text{proj}_v a)$$

$$H_v a = -(a - 2 \text{proj}_v a) = -\left(a - 2 \frac{v^T a}{v^T v} v\right)$$

$$= -\left(I_m - 2 \frac{vv^T}{v^T v}\right)a$$

How to use this for triangularization:



Step 1: Find constant c s.t.

$$H_v a = c \cdot e^1 \leftarrow \text{first unit vector } \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

$$-a + 2 \frac{v^T a}{v^T v} v = c e^1$$

or equiv. find a constant c s.t.

$$a - 2 \frac{v^T a}{v^T v} v = c e^1$$

$$\Rightarrow v = (a - c e^1) \frac{v^T v}{2 v^T a}$$

$\leadsto v$ must be parallel to $a - c e^1$

Note: Scaling of v doesn't effect the formula

Set $v = a - c e^1$

$$\frac{v^T v}{2 v^T a} = 1$$

$$\frac{\|a - c e^1\|_2^2}{2 \langle a - c e^1, a \rangle} = 1$$

$$\cancel{\|a\|_2^2} - 2c \langle a, e^1 \rangle + c^2 = 2 \cancel{\|a\|_2^2} - 2c \langle a, e^1 \rangle$$

$$c = \pm \|a\|_2$$

Choose c s.t. cancellation is avoided

If a is almost parallel to e^1
either $a + \|a\|_2 e^1$ or $a - \|a\|_2 e^1$

is small \Rightarrow would be dividing by $\|v\|_2$ [could yield instability in HH reflections]

$$v = \begin{cases} \frac{1}{2}(a - \|a\|_2 e^1) & \text{if } a_1 < 0 \\ \frac{1}{2}(a + \|a\|_2 e^1) & \text{if } a_1 > 0 \end{cases}$$

For a matrix $A \in \mathbb{R}^{m,n}$: HH reflection on the first column

$$H_{v^1} a^1 = c e^1$$

$$\Rightarrow H_{v^1} A = \begin{bmatrix} c & & \\ 0 & \square & \\ \vdots & & \\ 0 & & \end{bmatrix}$$

↑ first step to triangularization

j -th step of HH:

vector a^j is split into $\begin{bmatrix} a_1^j \\ a_2^j \end{bmatrix} \in \mathbb{R}^{m-j}$ # we want to keep

Find v^j s.t. $H_{v^j} a^j = \begin{bmatrix} \tilde{a}_1^j \\ 0 \\ \vdots \\ 0 \end{bmatrix} \in \mathbb{R}^j$ } $m-j$ zeros

Choose $v^j = \begin{bmatrix} 0 \\ a_2^j \end{bmatrix} - c_j e^j$ with $c_j = \pm \|a_2^j\|$
↑ choose s.t. no cancellation

Calculation shows: $\frac{\langle a^j, v^j \rangle}{\|v^j\|_2^2} = \frac{1}{2}$

$$H_{v^j} a^j = a^j - 2 \frac{\langle a^j, v^j \rangle}{\|v^j\|_2^2} v^j = a^j - v^j$$

↑ ↑
coincide on
indices $j+1$ to m

$$\Rightarrow H_{v^j} a^j = \begin{bmatrix} \tilde{a}_1^j \\ 0 \\ \vdots \\ 0 \end{bmatrix} \left. \begin{array}{l} \} \in \mathbb{R}^j \\ \} m-j \text{ zeros} \end{array} \right\}$$

Note: $H_{v^j} q^k$ with $k < j$ $q^k = H_{v^k} a^k$

$$H_{v^j} q^k = q^k - 2 \frac{\langle v^j, q^k \rangle}{\|v^j\|_2^2} v^j = q^k$$

$$\langle v^j, q^k \rangle = 0 \quad k \leq j-1$$

first $j-1$ entries zero only first k entries nonzero

Altogether:

$$H_{v^n} \dots H_{v^1} A = R$$

\ / / /
orthogonal

$$\Rightarrow Q = H_{v^1}^T \dots H_{v^n}^T \text{ orth.}$$

v^j : first $j-1$ entries are zero

Q is stored implicitly by storing vectors v^1, \dots, v^n as lower triangular matrix (compressed format)

Householder reflections in EIGEN:

C++-code 3.3.34: QR-decompositions in EIGEN

```

2 # include <Eigen/QR>
3
4 // Computation of full QR-decomposition (3.3.3.1),
5 // dense matrices built for both QR-factors (expensive!)
6 std::pair<MatrixXd, MatrixXd> qr_decomp_full(const MatrixXd& A) {
7   Eigen::HouseholderQR<MatrixXd> qr(A);
8   MatrixXd Q = qr.householderQ(); // apply n times HH reflections
9   MatrixXd R = qr.matrixQR().template triangularView<Eigen::Upper>();
10  return std::pair<MatrixXd, MatrixXd>(Q,R);
11 }
12
13 // Computation of economical QR-decomposition (3.3.3.1),
14 // dense matrix built for Q-factor (possibly expensive!)
15 std::pair<MatrixXd, MatrixXd> qr_decomp_eco(const MatrixXd& A) {
16   using index_t = MatrixXd::Index;
17   const index_t m = A.rows(), n = A.cols();
18   Eigen::HouseholderQR<MatrixXd> qr(A);
19   MatrixXd Q = (qr.householderQ() * MatrixXd::Identity(m,n)); // apply n times HH reflections
20   MatrixXd R = qr.matrixQR().block(0,0,n,n).template
21     triangularView<Eigen::Upper>(); //
22   return std::pair<MatrixXd, MatrixXd>(Q,R);
23 }

```

In general: only HH reflections are applied instead of building Q

Asymptotic complexity of Householder QR-decomposition

The computational effort for **HouseholderQR()** of $A \in \mathbb{R}^{m,n}$, $m > n$, is $O(mn^2)$ for $m, n \rightarrow \infty$.

each HH reflection applied to one vector $\sim O(m)$

applied to A $\sim O(mn)$

n such reflections \rightarrow overall $O(n \cdot mn)$

Normal equations vs. orthogonal transformations method

- Superior numerical stability (\rightarrow Def. 1.5.85) of orthogonal transformations methods:
- ▶ Use orthogonal transformations methods for least squares problems (3.1.38), whenever $A \in \mathbb{R}^{m,n}$ dense and n small.
- SVD/QR-factorization cannot exploit sparsity:
- ▶ Use normal equations in the expanded form (3.2.8)/(3.2.9), when $A \in \mathbb{R}^{m,n}$ sparse (\rightarrow Notion 2.7.1) and m, n big.

Alternative method for QR factorization:

Given rotations [build Q through rotations]

3.4. Singular Value Decomposition

A different orth. decomposition

Theorem 3.4.1. singular value decomposition → [?, Thm. 9.6], [?, Thm. 11.1]

For any $A \in \mathbb{K}^{m,n}$ there are unitary matrices $U \in \mathbb{K}^{m,m}$, $V \in \mathbb{K}^{n,n}$ and a (generalized) diagonal (*) matrix $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_p) \in \mathbb{R}^{m,n}$, $p := \min\{m, n\}$, $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_p \geq 0$ such that

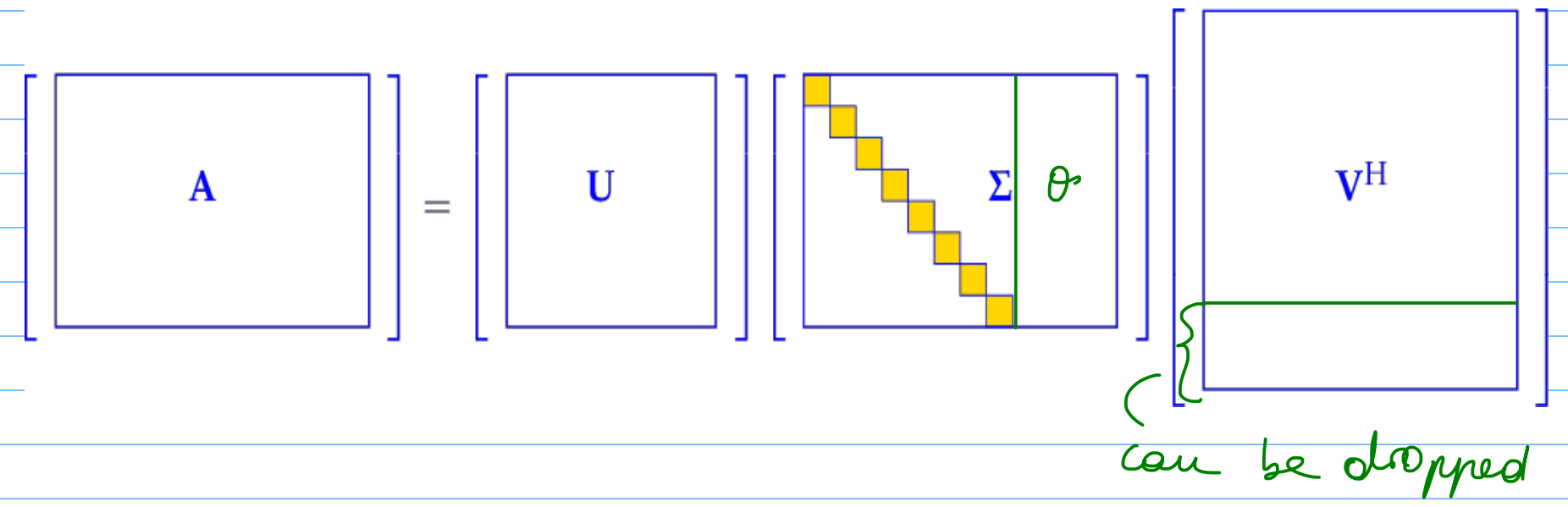
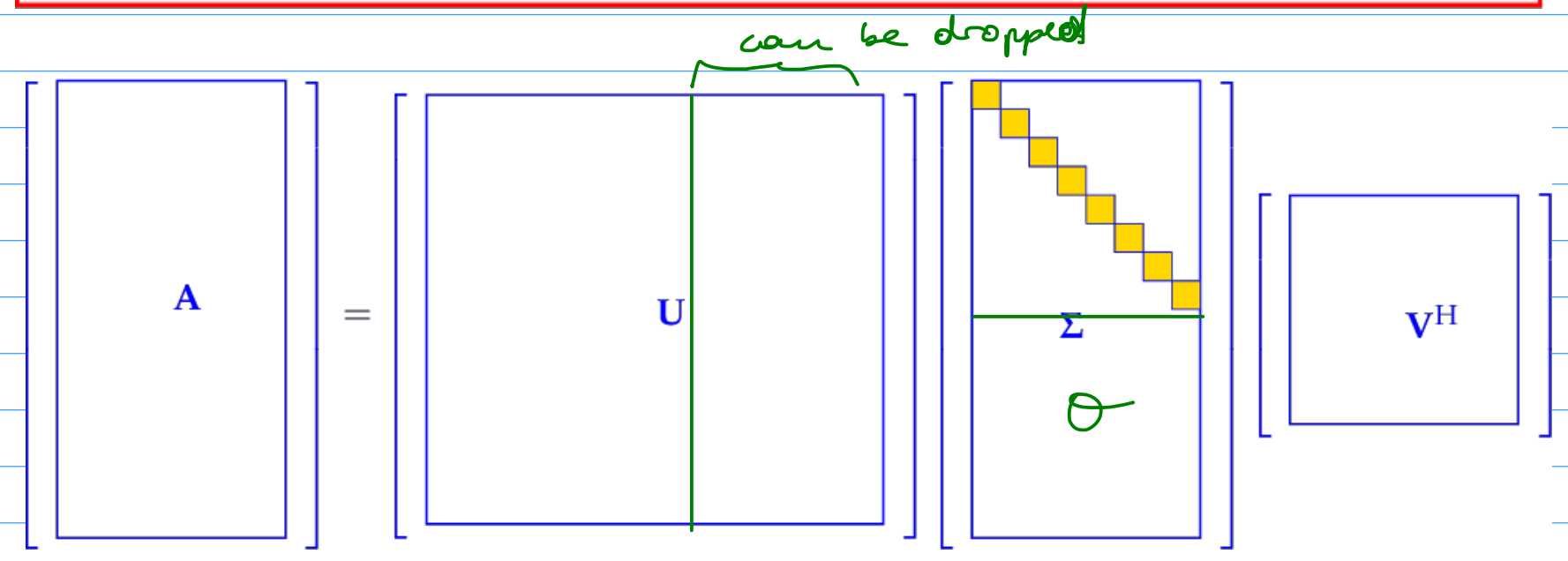
$A = U\Sigma V^H$.

$$U^H U = I_m = U U^H$$

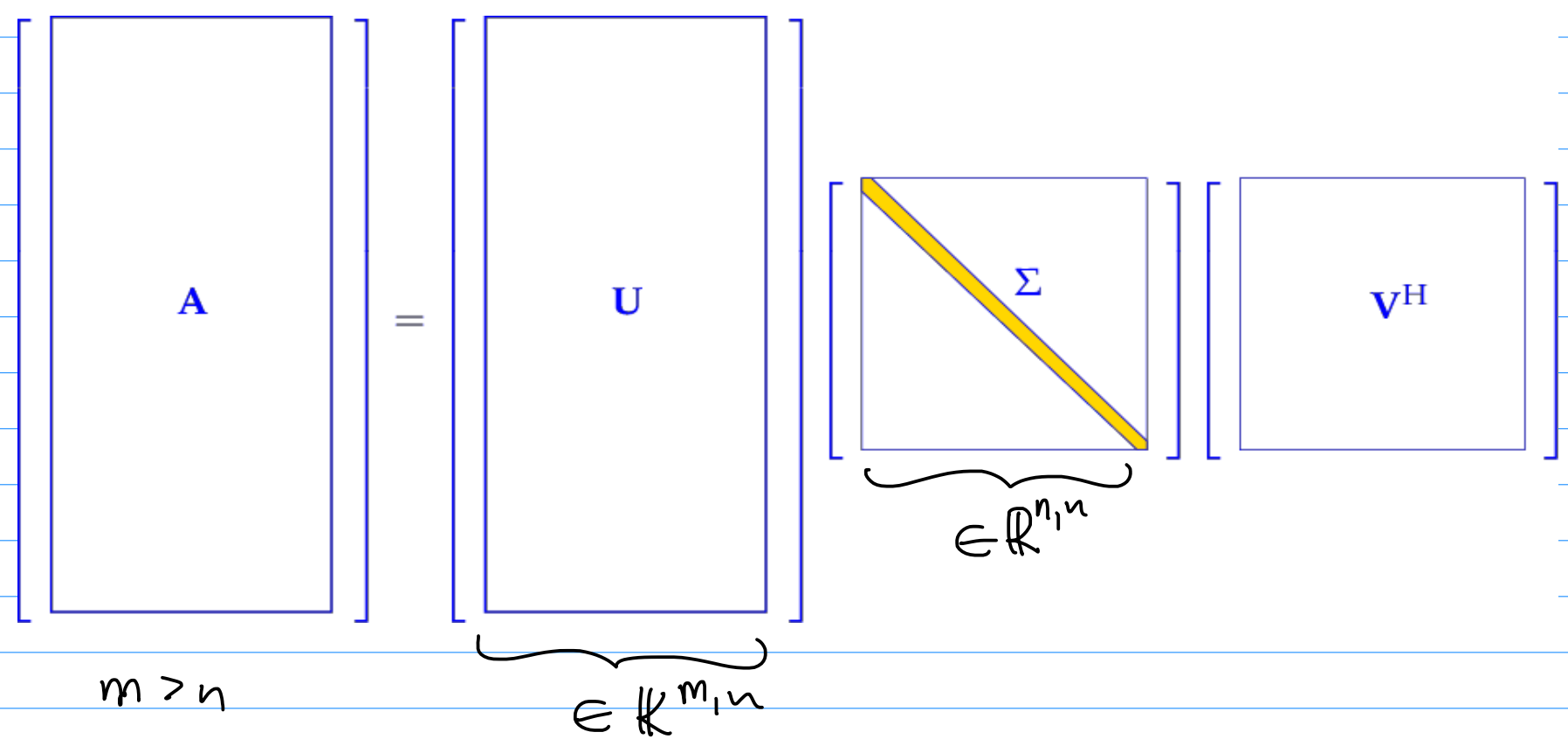
$$V^H V = I_n = V V^H$$

Definition 3.4.3. Singular value decomposition (SVD)

The decomposition $A = U\Sigma V^H$ of Thm. 3.4.1 is called **singular value decomposition (SVD)** of A . The diagonal entries σ_i of Σ are the **singular values** of A .



"Reduced / thin SVD" [as for QR]



$m < n$: $\Sigma_s \in \mathbb{R}^{m,m}$ $V \in \mathbb{K}^{n,m}$

Lemma 3.4.11. SVD and rank of a matrix → [?, Cor. 9.7]
 If, for some $1 \leq r \leq p := \min\{m, n\}$, the singular values of $A \in \mathbb{K}^{m,n}$ satisfy $\sigma_1 \geq \dots \geq \sigma_r > \sigma_{r+1} = \dots = \sigma_p = 0$, then

- $\text{rank}(A) = r$ (no. of non-zero singular values),
- $\mathcal{N}(A) = \text{Span}\{(\mathbf{V})_{:,r+1}, \dots, (\mathbf{V})_{:,n}\}$,
- $\mathcal{R}(A) = \text{Span}\{(\mathbf{U})_{:,1}, \dots, (\mathbf{U})_{:,r}\}$.

→ SVD reveals rank of A!
 → U, V encode orthonormal representations of $\mathcal{N}(A), \mathcal{R}(A)$

Lemma 3.4.5.
 The squares σ_i^2 of the non-zero singular values of A are the non-zero eigenvalues of $A^H A, A A^H$ with associated eigenvectors $(\mathbf{V})_{:,1}, \dots, (\mathbf{V})_{:,p}, (\mathbf{U})_{:,1}, \dots, (\mathbf{U})_{:,p}$, respectively.

$$\underbrace{A}_{\in \mathbb{R}^{m,n}} = \underbrace{[U_1 \quad U_2]}_{\in \mathbb{R}^{m,m}} \underbrace{\begin{bmatrix} \Sigma_r & 0 \\ 0 & 0 \end{bmatrix}}_{\in \mathbb{R}^{m,n}} \underbrace{\begin{bmatrix} V_1^H \\ V_2^H \end{bmatrix}}_{\in \mathbb{R}^{n,n}}$$

(3.4.22)

• $\text{span} \{ (V)_{:,r+1}, \dots, (V)_{:,n} \} = \mathcal{N}(A)$:

Take any $(V)_{:,j}$ $j \in \{r+1, \dots, n\}$

$$A (V)_{:,j} = U \Sigma \underbrace{V^T (V)_{:,j}}_{\uparrow \begin{bmatrix} \Sigma_r & 0 \\ 0 & 0 \end{bmatrix}} = U \Sigma e^j$$

$$\begin{matrix} j > r \\ \Rightarrow U \begin{bmatrix} \Sigma_r & 0 \\ 0 & 0 \end{bmatrix} e^j = 0 \end{matrix}$$

$$\Rightarrow \forall j \in \{r+1, \dots, n\}: A (V)_{:,j} = 0$$

$$\Rightarrow (V)_{:,j} \in \mathcal{N}(A)$$

$\{ (V)_{:,r+1}, \dots, (V)_{:,n} \}$ ONB for $\mathcal{N}(A)$

• columns of U_1 to span $\mathcal{R}(A)$

$$y \in \mathcal{R}(A) \Leftrightarrow \exists x \in \mathbb{R}^n \text{ s.t. } Ax = y$$

$$U \Sigma V^T x = y$$

$$\begin{bmatrix} u_1 & u_2 \end{bmatrix} \begin{bmatrix} \Sigma_r & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} v_1^T \\ v_2^T \end{bmatrix} x = y$$

$$u_1 \underbrace{\Sigma_r v_1^T}_{=: c \in \mathbb{R}^r} x = y \quad c = \begin{pmatrix} c_1 \\ \vdots \\ c_r \end{pmatrix}$$

$$\Rightarrow y = c_1 \underbrace{(u_1)_{:,1}}_{=(u)_{:,1}} + \dots + c_r \underbrace{(u_1)_{:,r}}_{=(u)_{:,r}}$$

for every $y \in \mathcal{R}(A)$:

$$\exists c \in \mathbb{R}^r \text{ s.t. } y = \sum_{j=1}^r c_j (u)_{:,j}$$

$$\Rightarrow \text{span} \{ (u)_{:,1}, \dots, (u)_{:,r} \} = \mathcal{R}(A)$$

3.4.2. SVD in EIGEN

JacobiSVD

C++-code 3.4.13: Computing SVDs in EIGEN

```

2 # include <Eigen/SVD>
3
4 // Computation of (full) SVD  $A = U\Sigma V^H \rightarrow$  Thm. 3.4.1
5 // SVD factors are returned as dense matrices in natural order
6 std::tuple<MatrixXd, MatrixXd, MatrixXd> svd_full(const MatrixXd& A) {
7     Eigen::JacobiSVD<MatrixXd> svd(A, Eigen::ComputeFullU |
8         Eigen::ComputeFullV);
9     MatrixXd U = svd.matrixU(); // get unitary (square) matrix U
10    MatrixXd V = svd.matrixV(); // get unitary (square) matrix V
11    VectorXd sv = svd.singularValues(); // get singular values as vector
12    MatrixXd Sigma = MatrixXd::Zero(A.rows(), A.cols());
13    const unsigned p = sv.size(); // no. of singular values
14    Sigma.block(0,0,p,p) = sv.asDiagonal(); // set diagonal block of  $\Sigma$ 
15    return std::tuple<MatrixXd, MatrixXd, MatrixXd>(U, Sigma, V);
16 }
17 // Computation of economical (thin) SVD  $A = U\Sigma V^H$ , see (3.4.4)
18 // SVD factors are returned as dense matrices in natural order
19 std::tuple<MatrixXd, MatrixXd, MatrixXd> svd_eco(const MatrixXd& A) {
20     Eigen::JacobiSVD<MatrixXd> svd(A, Eigen::ComputeThinU |
21         Eigen::ComputeThinV);
22     MatrixXd U = svd.matrixU(); // get unitary (square) matrix U
23     MatrixXd V = svd.matrixV(); // get unitary (square) matrix V
24     VectorXd sv = svd.singularValues(); // get singular values as vector
25     MatrixXd Sigma = sv.asDiagonal(); // build diagonal matrix  $\Sigma$ 
26     return std::tuple<MatrixXd, MatrixXd, MatrixXd>(U, Sigma, V);
27 }

```


without the compute U & V flags:

only singular values are computed

"Numerical rank" e.g. rank()

$$r := \# \{ \sigma_i : |\sigma_i| \geq \text{tol} \cdot \max_j \{ |\sigma_j| \} \}$$

Default: tol = EPS setThreshold()

to set tol manually

Cost of thin SVD: $\mathcal{O}(mn^2)$ $m > n$

Jacobi SVD is numerically stable

3.4.3 Generalized solutions by SVD

Assume $A \in \mathbb{K}^{m,n}$ $m \geq n$ $\text{rank}(A) = r \leq n$

As before:

$$A = [u_1 \ u_2] \begin{bmatrix} \Sigma_r & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} V_1^H \\ V_2^H \end{bmatrix}$$

least squares: $\min \|Ax - b\|_2$

$$\|Ax - b\|_2^2 = \|U \Sigma V^H x - b\|_2^2 = \|U^H (U \Sigma V^H x - b)\|_2^2$$

$$= \|\Sigma V^H x - U^H b\|_2^2$$

$$= \left\| \begin{bmatrix} \Sigma_r V_1^H x \\ 0 \end{bmatrix} - \begin{bmatrix} u_1^H b \\ u_2^H b \end{bmatrix} \right\|_2^2$$

$$= \left\| \begin{bmatrix} \Sigma_r V_1^H x - U_1^H b \\ -U_2^H b \end{bmatrix} \right\|_2^2 = \left\| \Sigma_r V_1^H x - U_1^H b \right\|_2^2 + \underbrace{\|U_2^H b\|_2^2}_{\text{fixed}}$$

\uparrow $\|\cdot\|_{\mathbb{K}^m}$ $\|\cdot\|_{\mathbb{K}^r}$ $\|\cdot\|_{\mathbb{K}^{m-r}}$

equivalently: minimize $\left\| \Sigma_r V_1^H x - U_1^H b \right\|_2^2$

choose x s.t.

$$\Sigma_r V_1^H x = U_1^H b \quad \text{LSE}$$

$r \times n$
(possibly underdetermined)

If $r < n$: generalized solution

\leadsto non-uniqueness from

$$\mathcal{N}(V_1^H)$$

$$\mathcal{N}(V_1^H) = \mathcal{Q}(V_1)^\perp = \mathcal{Q}(V_2)$$

choose $x \in \mathcal{Q}(V_1)$

$$x = V_1 z \quad \text{for some } z \in \mathbb{K}^r$$

$$\Sigma_r \underbrace{V_1^H V_1}_{=I} z = U_1^H b$$

$$z = \Sigma_r^{-1} U_1^H b$$

Generalized solution: $x = V_1 z = V_1 \Sigma_r^{-1} U_1^H b$

Theorem 3.4.30. Pseudoinverse and SVD

If $A \in \mathbb{K}^{m,n}$ has the SVD decomposition $A = U \Sigma V^H$ partitioned as in (3.4.22), then its Moore-Penrose pseudoinverse (\rightarrow Thm. 3.1.37) is given by $A^\dagger = V_1 \Sigma_r^{-1} U_1^H$.

[this requires $\text{rank}(A) \rightarrow$ determine numerical rank]