# Numerical Methods for CSE

### Autumn 2023
### Q+A

22.09.2023

VIDEOS

LECTURE NOTES

↓ ↘
HAND SCRIPT

↓

every week.

interactive

Q+A

↓

discona

HOMEWORK *

TUTORIALS.

very important

very important

NO BONUS

A MOCK EXAM   "BYOD" + "DIY"   NOVEMBER

"CODE EXPERT"

# Def kronecker produkt of two matrices

$A \in \mathbb{R}^{m \times n}$  $B \in \mathbb{R}^{l \times k}$

$m, n, l, k \in \mathbb{N}$

$A \otimes B \in \mathbb{R}^{(ml) \times (nk)}$   Block of size $l \times k$

$$\begin{bmatrix} A_{11}B & A_{12}B & \cdots & A_{1n}B \\ \vdots & & & \\ A_{m1}B & A_{m2}B & \cdots & A_{mn}B \end{bmatrix}$$

$(ml) \times (nk)$

**Example**

$$I_m \otimes A = \begin{bmatrix} 1A & & 0 \\ & 1A & \\ 0 & & \ddots \\ & & & 1A \end{bmatrix}$$

$$A \times I_m = \begin{bmatrix} a_{11}I & a_{12}I & \cdots & a_{1k} \\ & a_{22}I & & \\ \vdots & & & \\ a_{n1} & a_{n2} & \cdots & a_{nk}I \end{bmatrix}$$

$A_{n \times k}$

$$\begin{bmatrix} a_{11} & 0 & \cdots & 0 \\ 0 & a_{11} & & \vdots \\ \vdots & & \ddots & 0 \\ 0 & \cdots & & a_{11} \end{bmatrix} \begin{bmatrix} a_{12} & 0 & \cdots & 0 \\ 0 & a_{12} & & 0 \\ \vdots & & \ddots & 0 \\ 0 & \cdots & & a_{12} \end{bmatrix}$$

# Partial Differential Equations



$$\boxed{\frac{\partial}{\partial x}} \quad \frac{\partial}{\partial y}$$

$\downarrow$

Matrix

$$\frac{\partial}{\partial x} \quad \frac{\partial}{\partial y}$$

$$\underline{\underline{D}}_x \otimes \underline{\underline{D}}_y$$



$x_{i-1} \quad x_i \quad x_{i+1}$

$$\underline{u} = \begin{bmatrix} u(x_0) \\ \vdots \\ u(x_{N-1}) \end{bmatrix} \qquad \underline{\underline{A}} \quad \underline{\underline{D}}_x$$

$$\underline{\underline{D}}_y$$

$$\frac{\partial}{\partial x} u(x_i) \approx \frac{u_{i+1} - u_{i-1}}{2h}$$

# Q3 0.1.11.C.

**(Q3.0.1.11.C)** [A $\int_0^1 e^x\, de$-type problem]    We know the solution $\mathbf{x} \in \mathbb{R}^n$ and the right-hand-side vector $\mathbf{b} \in \mathbb{R}^n$ of the $n \times n$ (Toeplitz) tridiagonal linear system of equations

$$
\begin{bmatrix}
\alpha & \beta & 0 & \cdots & & & \cdots & 0 \\
\beta & \alpha & \ddots & & & & & \vdots \\
0 & \beta & \ddots & \ddots & & & & \\
\vdots & & \ddots & \ddots & \ddots & & & \\
 & & & \ddots & \ddots & \ddots & & \\
 & & & & \ddots & \ddots & \beta & 0 \\
\vdots & & & & & \beta & \alpha & \beta \\
0 & \cdots & & & & \cdots & 0 & \beta & \alpha
\end{bmatrix} \mathbf{x} = \mathbf{b}.
$$

Which overdetermined linear system of equations of maximal size has the vector $[\alpha, \beta]^\top \in \mathbb{R}^2$ as its solution?

---

$$\alpha x_1 + \beta x_2 = b_1$$
$$\beta x_1 + \alpha x_2 + \beta x_3 = b_2$$
$$\beta x_2 + \alpha x_3 + \beta x_4 = b_3$$
$$\beta x_3 + \alpha x_4 + \beta x_5 = b_4$$
$$\beta x_4 + \alpha x_5 + \beta x_6 = b_5$$
$$\vdots$$
$$\beta x_{n-2} + \alpha x_{n-1} + \beta x_n = b_{n-1}$$
$$\beta x_{n-1} + \alpha x_n = b_n$$

$$
\begin{bmatrix}
x_1 & x_2 \\
x_2 & x_1 + x_3 \\
x_3 & x_2 + x_4 \\
x_4 & x_3 + x_5 \\
\vdots & \vdots \\
x_{n-1} & x_{n-2} + x_n \\
x_n & x_{n-1}
\end{bmatrix}
\begin{bmatrix} \alpha \\ \beta \end{bmatrix}
=
\begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ \vdots \\ b_n \end{bmatrix}
$$

29.09.2023

**(Q3.1.1.14.C)** Given a matrix $\mathbf{B} \in \mathbb{R}^{m,n}$, a vector $\mathbf{c} \in \mathbb{R}^m$, and $\lambda > 0$, define

$$\{\mathbf{x}^*\} := \operatorname*{argmin}_{\mathbf{x} \in \mathbb{R}^n} \|\mathbf{B}\mathbf{x} - \mathbf{b}\|_2^2 + \lambda\|\mathbf{x}\|_2^2 \subset \mathbb{R}^n .$$

looks like a penalisation :

$x \neq 0 \Rightarrow \lambda\|x\|_{L}^2$ clear increase !

but actually no deeper meaning than for trick

State an overdetermined linear system of equations $\mathbf{A}\mathbf{x} = \mathbf{b}$, of which $\mathbf{x}^*$ is a least-squares solution.

$\hookrightarrow$ a method to address linear least squares problem for $\underline{\underline{B}}$ with $\underline{\underline{B}}$ not full rank

(i.e. ~~(some of)~~ the columns of $\underline{\underline{B}}$ are linear dependent)

Pick a small $\lambda > 0$

$$\underline{\underline{A}} = \begin{bmatrix} \underline{\underline{B}} \\ \sqrt{\lambda}\, \underline{\underline{I}}_n \end{bmatrix} \qquad \underline{b} = \begin{bmatrix} \underline{c} \\ \underline{0} \end{bmatrix} \in \mathbb{R}^{m+n}$$

$\hookrightarrow$ linear indepnt columns of $\underline{\underline{A}}$

$$\operatorname*{argmin}_{\underline{x} \in \mathbb{R}^n} \| \underline{\underline{A}}\,\underline{x} - \underline{b} \|_2^2$$

$\lambda\|\underline{x}\|_{L}^2 = \langle \sqrt{\lambda}\,\underline{\underline{I}}\,\underline{x},\ \sqrt{\lambda}\,\underline{\underline{I}}\,\underline{x} \rangle$

$(m+n) \times n$   has Rank $n$.

$\Rightarrow$ QR-decomposition will work. !

Possible advantage:
if $\underline{\underline{B}}$ is sparse,
so is $\underline{\underline{A}}$, so QR-dec
via Givens-Rotation
might be much
less expensive
than SVD($\underline{\underline{B}}$).

Question: difference in advantages CSC/CRS

CSS: good for slicing columns

CRS: good for slicing rows

both are good for internal +, * (pointwise)

→ Matrix * vektor (might be faster)

Remark other formats are better for a fast construction.

---

**(Q2.7.1.5.E)** For a given matrix $A \in \mathbb{R}^{m,n}$, $m, n \in \mathbb{N}$, we define the square matrix

$$W_A := \begin{bmatrix} O_{m,m} & A \\ A^\top & O_{n,n} \end{bmatrix} \in \mathbb{R}^{m+n,m+n}.$$

Outline the implementation of an efficient C++ function

```cpp
void crsAtoW(std::vector<double> &val,
             std::vector<unsigned int> &col_ind,
             std::vector<unsigned int> &row_ptr);
```

whose arguments supply the three vectors defining the matrix $A$ in CRS format and which overwrites them with the corresponding vectors of the CRS-format description of $W_A$.

Remember that the CRS format of a matrix $A \in \mathbb{R}^{m,n}$ is defined by

$$\text{val}[k] = (A)_{i,j} \quad \Leftrightarrow \quad \begin{cases} \text{col\_ind}[k] = j, \\ \text{row\_ptr}[i] \leq k < \text{row\_ptr}[i+1], \end{cases} \quad 1 \leq k \leq \text{nnz}(A).$$

It may be convenient to use **std::vector::resize(n)** that resizes a vector so that it contains n elements. If n is smaller than the current container size, the content is reduced to its first n elements, removing those beyond (and destroying them). If n is greater than the current container size, the content is expanded by inserting at the end as many elements as needed to reach a size of n using their default value.

---

Most important: how to implement $\underline{A^\top}$ ?

This function does it:

```cpp
CRSMatrix sparse_transpose(const CRSMatrix& input) {
    CRSMatrix res{
        input.m,
        input.n,
        input.nz,
        std::vector<double>(input.nz, 0.0),
        std::vector<int>(input.nz, 0),
        std::vector<int>(input.m + 2, 0) // one extra
    };

    // count per column
    for (int i = 0; i < input.nz; ++i) {
        ++res.rowPtr[input.colIndex[i] + 2];
    }

    // from count per column generate new rowPtr (but shifted)
    for (int i = 2; i < res.rowPtr.size(); ++i) {
        // create incremental sum
        res.rowPtr[i] += res.rowPtr[i - 1];
    }

    // perform the main part
    for (int i = 0; i < input.n; ++i) {
        for (int j = input.rowPtr[i]; j < input.rowPtr[i + 1]; ++j) {
            // calculate index to transposed matrix at which we should p
            const int new_index = res.rowPtr[input.colIndex[j] + 1]++;
            res.val[new_index] = input.val[j];
            res.colIndex[new_index] = i;
        }
    }
    res.rowPtr.pop_back(); // pop that one extra

    return res;
}
```

**(Q2.6.0.25.F)** [Loss of stability] By direct block-wise Gaussian elimination we found the following solution formulas for a block-partitioned linear system of equations with $\mathbf{D} \in \mathbb{R}^{n,n}$, $\mathbf{c}, \mathbf{b} \in \mathbb{R}^n$, $\alpha \in \mathbb{R}$, $\mathbf{y} \in \mathbb{R}^{n+1}$:

$$\mathbf{Ax} = \begin{bmatrix} \mathbf{D} & \mathbf{c} \\ \mathbf{b}^\top & \alpha \end{bmatrix} \begin{bmatrix} \mathbf{x}_1 \\ \xi \end{bmatrix} = \mathbf{y} := \begin{bmatrix} \mathbf{y}_1 \\ \eta \end{bmatrix}, \tag{2.6.0.7}$$

$$\blacktriangleright \boxed{\xi = \frac{\eta - \mathbf{b}^T \mathbf{D}^{-1} \mathbf{y}_1}{\alpha - \mathbf{b}^\top \mathbf{D}^{-1} \mathbf{c}}}, \tag{2.6.0.8}$$

$$\mathbf{x}_1 = \mathbf{D}^{-1}(\mathbf{y}_1 - \xi \mathbf{c}).$$

Use these formulas to compute the solution of the $2 \times 2$ linear system of equations
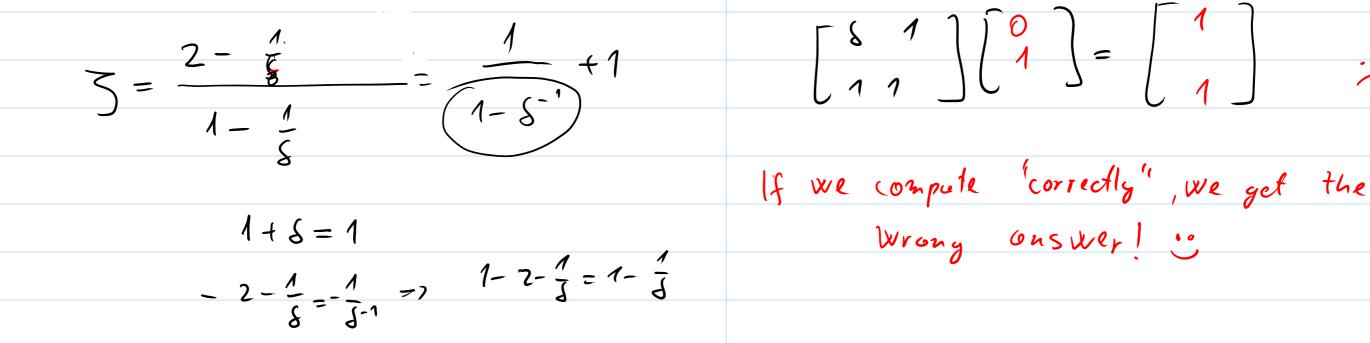
$$\begin{bmatrix} \delta & 1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ \xi \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \end{bmatrix},$$

assuming $|\delta| < \frac{1}{2}\text{EPS}$ and using floating point arithmetic.

**Hint.** Remember that, if $|\delta| < \frac{1}{2}\text{EPS}$, in floating point arithmetic

$$1 \dot{\mp} \delta = 1 \quad \text{and} \quad 2 \dot{\mp} \delta^{-1} = \delta^{-1}.$$

This is compatible with the "Axiom" of roundoff ana $\therefore$ Ass. 1.5.3.11

$$2 - \delta^{-1} = 2 - (2 + \delta^{-1}) = -\delta^{-1}$$
$$=$$

$$1 - \delta^{-1} = 1 - (2 + \delta^{-1})$$
$$= -1 - \delta^{-1}$$

$$\xi = \frac{-\delta^{-1}}{-\delta^{-1} - 1} = \frac{+1}{+1 + \delta^{-1}} = \frac{1}{1} = 1$$

$$x_1 = \delta^{-1}(1 - \quad 1) = 0 \implies \begin{bmatrix} -1 \\ 0 \end{bmatrix}$$

$$\xi = \frac{2 - \frac{1}{\xi}}{1 - \frac{1}{\delta}} = \frac{\frac{1}{\boxed{1 - \delta^{-1}}}}{} + 1$$

$$1 + \delta = 1$$

$$-2 - \frac{1}{\delta} = -\frac{1}{\delta^{-1}} \implies \quad 1 - 2 - \frac{1}{\delta} = 1 - \frac{1}{\delta}$$

$$\begin{bmatrix} \delta & 1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \because$$

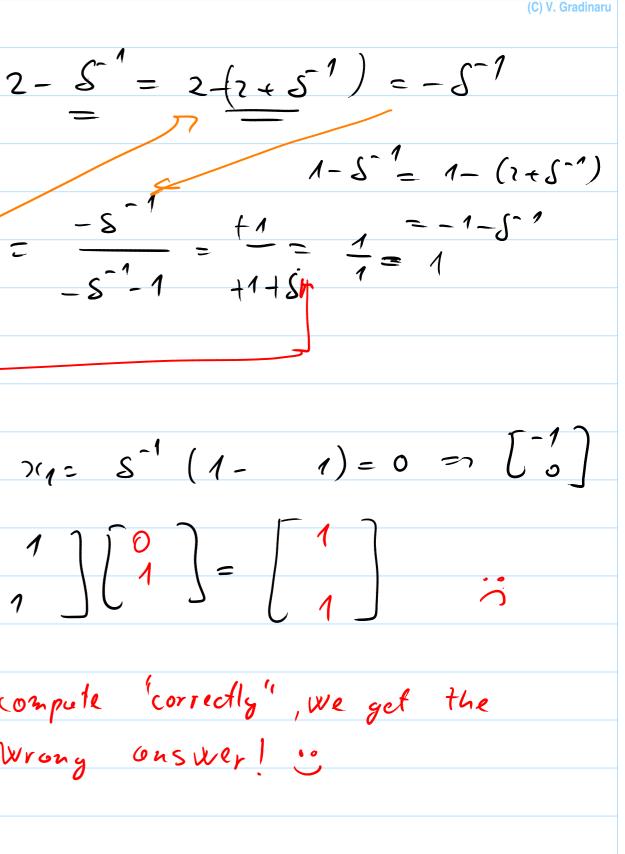If we compute 'correctly', we get the wrong answer! $\because$

Could you please explain what the algorithm 2.10.1 from the exercise sheet does?
I don't quite understand why we can use the Sherman-Morris-Woodbury Formula
for sub-problem c)

affine space $\neq$ linear space
vektor

$A_{,+,\cdot_{sc}}$      $V_{,+,\cdot_{sc.}}$
    $\ni 0$

$V = \text{Ker}(\underline{\underline{M}}) = \{\underline{x} \in \mathbb{R}^n ; \underline{\underline{M}}\underline{x} = \underline{0}\}$
    $\ni 0$

$A = \{\underline{x} \in \mathbb{R}^n ; \underline{\underline{M}}\underline{x} = \underline{b}\}$

$\underline{b}$

$\underline{x} \in A \iff \underline{x} = \underline{b} + \underline{v} \text{ with } \underline{v} \in V$

Use SMW-Formula in order
to avoid solving.

$\underline{\underline{M}}\underline{x} = \underline{ev}$ with a full

matrix $\underline{\underline{M}}$

Use the fact that $\underline{\underline{M}} = \text{diag}(\underline{d}) + \underline{u}\,\underline{v}^T$
    $\uparrow$

$\searrow$ Rang 1-modification

Sparse

$\text{diag}(\underline{d})\underline{x} = \underline{b}$ is solvable
directly in $O(n)$
Operations

Hence $O(n)$ instead of $O(n^3)$.

**(Q3.4.4.13.A)** Let $\mathbf{M} \in \mathbb{R}^{n,n}$ be symmetric and positive definite (*s.p.d.*) and $\mathbf{A} \in \mathbb{R}^{m,n}$. Devise an algorithm for computing

$$\operatorname*{argmax}_{\mathbf{x} \in B} \|\mathbf{Ax}\| \quad , \quad B := \{\mathbf{x} \in \mathbb{R}^n, \mathbf{x}^\top \mathbf{Mx} = 1\},$$

also based on the SVD of $\mathbf{M}$.

$$\underline{\underline{M}} = \underline{\underline{U}} \underline{\underline{\Sigma}} \underline{\underline{U}}^T \qquad \lambda_1 \geqslant \lambda_2 \geqslant \ldots \geqslant \lambda_n > 0$$

$$\underline{\underline{D}} = \sqrt{\underline{\underline{\Sigma}}}$$

Write the condition: $\underline{x}^T \underline{\underline{M}} \underline{x} = 1$

$$\underline{x}^T \underline{\underline{U}} \underline{\underline{\Sigma}} \underline{\underline{U}}^T \underline{x} = 1$$

$$\underbrace{\underline{x}^T \underline{\underline{U}} \underline{\underline{D}}^T}_{\underline{y}^T} \underbrace{\underline{\underline{D}} \underline{\underline{U}}^T \underline{x}}_{\underline{y}} = 1$$

orthonormal.
↓

$$\underline{y}^T \underline{y} = 1 \quad \text{with} \quad \underline{y} = \underline{\underline{D}} \underline{\underline{U}}^T \underline{x} \iff \underline{\underline{U}} \underline{\underline{D}}^{-1} \underline{y} = \underline{x}$$

$$\underline{\underline{D}}^{-1} |$$

$$\underline{\underline{U}} |$$

$$\operatorname*{argmax}_{\underline{x} \in \mathbb{R}^n} \|\underline{\underline{A}} \underline{x}\| =$$

$$\underline{x}^T \underline{\underline{M}} \underline{x} = 1$$

$$\operatorname{argmax} \|\underline{\underline{A}} \underline{\underline{U}} \underline{\underline{D}}^{-1} \underline{y}\| = \underline{\underline{U}} \underline{\underline{D}}^{-1} \operatorname*{argmax}_{\underline{y} \in \mathbb{R}^n} \|\underline{\underline{B}} \underline{y}\|$$

$$\underline{y} = \underline{\underline{D}} \underline{\underline{U}}^T \underline{x}$$
$$\underline{y}^T \underline{y} = 1 \qquad \underbrace{\qquad}_{\underline{\underline{B}}} \qquad \underline{y}^T \underline{y} = 1$$

Which $\underline{x}$ makes $\|\underline{\underline{A}} \underline{x}\| = \min$

See (3.4.4.3) for the solution of this problem.

Sol. is the first right singular vector of $\underline{\underline{B}}$

Sol $\underline{x} = \underline{\underline{U}} \underline{\underline{D}}^{-1} \underline{y}$

**Q)** Why $\det H = -1$ for any Householder Matrix?



$\underline{n} \in \mathbb{R}^d$

$\|\underline{n}\| = 1$

$\|\underline{n}\| = 1$

plane $\perp \underline{n}$

$\dim V = d-1$

Take basis (orthonormal) $\underline{v_1}, \cdots, \underline{v_{d-1}} \in V$

$$\underline{\underline{H}}\, \underline{v_1} = \underline{v_1} \;,\; \ldots \;,\; \underline{\underline{H}}\, \underline{v_{d-1}} = \underline{v_{d-1}}$$

$$\underline{\underline{H}}\, \underline{n} = -\underline{n}$$

$\rightarrow$ 1 is EW of $\underline{\underline{H}}$ of Multiplicity $d-1$

$-1$ is EW of $\underline{\underline{H}}$ of Multiplicity 1

$$\det \underline{\underline{H}} = \lambda_1 \lambda_2 \lambda_2 \cdots \lambda_n = (1)^{d-1}(-1) = -1$$

---

**Q)** How do we solve Linear Systems of Eq.?

1) Want good precision? $\neq$ Moderate precision.

2) how expensive?

good precision, not too big Matrix $\Rightarrow$ LU-decomp.

A symmetric pos.df $\Rightarrow$ Cholesky-decomp

big Matrix, sparse (banded)

$\rightarrow$ there might be some direct methods that keep sparsity, so might be feasable $\quad$ QR; LU

in general LU would be not feasable.

less precision or Large sparse Matrix $\underline{\underline{A}}$

$\Rightarrow$ iterative methods Krylov-type Method, $\begin{cases} \text{symmetric} & CG \\ \text{non-symm} & \text{others} \end{cases}$

Krylov-type methods use only

$$\underset{A}{\underline{A}} \underset{x}{\underline{x}}$$

$\Rightarrow$ if $\underline{A}$ is sparse $\Rightarrow$ $O(2, n^2)$

Note : CG slow if $\text{cond}(\underline{A})$ is big.

$\rightarrow$ use pre-conditioning

$$\underline{A}\underline{x} = \underline{b} \Rightarrow \underline{B}^{-1}\underline{A}\underline{x} = \underline{B}^{-1}\underline{b}$$

$$\underline{B}^{-1} \mid$$

if $\underline{B}^{-1} = \underline{A}^{-1} \Rightarrow$ $\underline{x} = \underline{B}^{-1}\underline{b}$ ☺

use $\underline{B}^{-1} \approx \underline{A}^{-1}$

ex $\underline{B} = \text{diag}(\underline{A}) \rightarrow \underline{B}^{-1}$ cheap

$\Rightarrow$ solve (CG) for $\underline{B}^{-1}\underline{A}\underline{x} = \underline{B}^{-1}\underline{b}$

---

Ⓠ Polar decomposition

$$\underline{A} = U\Sigma V^T = \underbrace{UV^T}_{\underline{Q}}\underbrace{V\Sigma V^T}_{H \text{ spd.}}$$

$$\underset{\wedge}{} \quad \mathbb{1} = V^T V$$

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

$$\underbrace{\qquad}_{\tilde{H} \text{ spd.}} \underbrace{\qquad}_{\tilde{Q}}$$

$m \geq n$ economical SVD

$$U \in \mathbb{R}^{m \times n}, \quad \Sigma = \begin{bmatrix} \sigma_1 & 0 \\ & \ddots & \\ 0 & \sigma_n \end{bmatrix}$$

$$V \in \mathbb{R}^{n \times n}$$

$$V\Sigma V^T \sim n^2 \text{ Operations (storage)}$$

$$U \sim m \cdot n \text{ Operations (storage)}$$

$\Rightarrow$ not cheaper $O(\widehat{mn} + n^2)$

$+ n \rightarrow$ dominates !

3.12.d. Main messages:

① avoid using economical QR-decomposition

② if using it, then be aware of dimensions!

3.$\pi$.  $\quad X_1 = A_1 B_1^T \in \mathbb{R}^{m\times n} \qquad X_2 = A_2 B_2^T$

$\qquad\qquad\qquad\qquad\qquad\qquad A_1, A_2 \in \mathbb{R}^{m\times k}$

$X = [X_1 \; X_2] \in \mathbb{R}^{m,2n} \qquad\qquad B_1, B_2 \in \mathbb{R}^{n\times k}$

$R(X) = R(X_1) + R(X_2)$

$\Rightarrow \dim R(X) \leq \dim R(X_1) + \dim R(X_2) = 2k \Rightarrow$
$\qquad\qquad\qquad\qquad \, {}^{\leq k} \qquad\qquad {}^{\leq k}.$

$\dim R(X) \leq m$

$\dim R\left(\big|\big|\big|\right) = 1 \qquad \dim R\left(\big|\big|\big| \;\equiv\right) = k = \text{Rank}(A)$

$\dim R(X) \leq \min\{m, 2k\}.$

SVD of $\underline{\underline{X}}$ without assembling X?
$\qquad\qquad$ i.e. by using only the factors
$\qquad\qquad\qquad A_1, B_1, A_2, B_2.$

economical!

$B_1 = Q_1 R_1$

$B_2 = Q_2 R_2$

$X = [A_1 R_1^T Q_1^T \quad A_2 R_2^T Q_2^T] =$

$= \underbrace{[A_1 R_1^T \quad A_2 R_2^T]}_{\Sigma \; {}^{m\times 2n}} \underbrace{\begin{bmatrix} Q_1^T & 0 \\ 0 & Q_2^T \end{bmatrix}}_{Q \in \mathbb{R}^{2k\times 2n}}$

${}^{m\times n} \qquad {}^{m\times n}$

$\Sigma = \tilde{U} \tilde{\Sigma} \tilde{V}^T \qquad \tilde{U} \in \mathbb{R}^{m,2k}, \; \tilde{V} \in \mathbb{R}^{2k\times 2n}$

$$X = \tilde{U} \, \tilde{\Sigma} \, \tilde{V}^T \, Q^T$$

$$U'' \quad \Sigma' \quad Q$$

**Implementation:**

$$I_{n,k} = \begin{bmatrix} 1 \\ & 1 \\ & & \ddots \\ & & & 1 \end{bmatrix}$$

i.s.b.

one call of full QR-decomposition (Householder)
gives whole $n \times n$ matrix $\underline{Q}_1$

$$\underline{Q}_1 \, \underline{I}_{n,k} = \begin{bmatrix} & & \\ & & \\ & & \end{bmatrix} = \text{first } k \text{ columns of } \underline{Q}.$$

$$\underline{q}_1 \cdots \underline{q}_k$$

---

(Q) Meaning of interpolation operator?

$t_1 \quad \underline{y}_1 \in \mathbb{R}^n$

$t_2 \quad \underline{y}_2 \in \mathbb{R}^n$

$\vdots$

$t_m \quad \underline{y}_m$



Want a function $f : [t_1, t_m] \to \mathbb{R}^n$ such that

$$f(t_j) = \underline{y}_j \quad \text{for all } j = 1, \dots, m.$$

$f$ should have some desired properties

applications          mathematical feasable

$$V = \text{span} \{ b_1, b_2, \dots, b_n \} \ni f$$

$$f(t) = \sum_{k=1}^{n} c_k \, b_k(t)$$

$n$ conditions $\Rightarrow$ Linear System $n \times n$ for $c_1, \dots, c_n$.

Which $b_1, \ldots b_n$ to choose?

$V = P_n \implies b_k(t) = t^{k-1}$ ∴

↳

linear system bad conditioned!

+ better basis.!

+ suitable space.

$b_k(t) = e^{ikt} = \cos(kt) + i \sin(kt)$

$\implies$ trigonometric polynomials

$\implies$ very fast and accurate algorithms

via FFT.

↳ complete basis in $L^2$

---

Ⓠ Lagrange interpolation.

Let us take 3 measurement.

| $t_0$ | $t_1$ | $t_2$ |
|-------|-------|-------|
| $y_0$ | $y_1$ | $y_2$ |



$$L_0(t) = \frac{t-t_1}{t_0-t_1} \frac{t-t_2}{t_0-t_2} \implies L_0(t_0) = \boxed{\frac{t_0-t_1}{t_0-t_1}} \boxed{\frac{t_0-t_2}{t_0-t_2}} = 1$$

$$L_0(t_1) = \frac{\textcolor{red}{0}}{t_0-t_1} \frac{t_1-t_2}{t_0-t_2} \implies L_0(t_1) = L_0(t_2) = 0$$



Parabola since polynom. of degree 2.

(C) V. Gradinaru

$f(t) = y_0 L_0(t) + y_1 L_1(t) + y_2 L_2(t)$

$f(t_0) = y_0 \cdot 1 + y_1 \cdot 0 + y_2 \cdot 0 = y_0$

$f(t_1) = y_0 \cdot 0 + y_1 \cdot 1 + y_2 \cdot 0 = y_1$

$f(t_2) = y_0 \cdot 0 + y_1 \cdot 0 + y_2 \cdot 1 = y_2$

$\uparrow$

"$\delta$"-Property.    $L_j(t_k) = \begin{cases} 1 & \text{if } j=k \\ 0 & \text{if } j \neq k. \end{cases}$