

# 1) Errors and numbers

## 1.1 Preliminary considerations

There are three main sources of error in the numerical analysis of physical systems:

- i) rounding errors, since instead of numbers  $t \in \mathbb{R}$  a finite set of numbers is calculated.
- ii) "discretization" or "truncation error"

Examples:

- $\int_a^b f(x) dx$  is replaced by a finite sum ( $dx \rightarrow \Delta x$ )
- Power series  $P(x) = \sum_{n=1}^{\infty} P_n x^n$  is approximated by a finite sum

- iii) physical approximation errors:

usually physical approximations are made before the numerical evaluation. Examples from solid state physics: Hartree-Fock approximation, density functional theory, ...

Estimation of the order of magnitude of the errors:

- iii) usually difficult
- ii) order of magnitude can often be estimated. If you know, for example, that a property  $F(\Delta x)$  is linear in  $\Delta x$  for  $\Delta x \rightarrow 0$  then one may calculate  $F_i = F(\Delta x_i)$  for some values  $\Delta x_i$ :
  - $\Rightarrow$  linear fit  $F^L(\Delta x) = F(0) + m \Delta x$
  - $\Rightarrow F(0)$  the property one is looking for plus statistical error
- i) We consider these errors in this chapter

## 1.2 Rounding errors

### 1.2.1 Representation of numbers

A) Integer numbers

Computers construct whole ("integer") numbers in the form

$$x = \pm (2^m 2^n + 2^{m-1} 2^{n-1} + \dots + 2^0 2^0) = \pm \sum_{i=0}^m a_i 2^i$$

with  $a_i = 0$  or  $a_i = 1$  (value of bit  $i$ )

With 16 bits: Bit 1 determines  $\pm$

Bit 2-16 determines  $|x|$

$\Rightarrow \exists 2^{15}$  such numbers between 0 and  $2^{15}-1 = 32767$

so  $x \in \{-32767, \dots, -1, 0, 1, \dots, 32767\}$

Modern programming languages translate this coding into the decimal system

$$x = \pm \sum_{i=0}^m a_i \cdot 10^i$$

For the sake of simplicity, we will only consider the decimal system in the following. However, all considerations apply completely analogously to the binary system.

B) "floating point" = "fp" representation

Question: How do you represent real numbers (approximately)?

Idea: by using rational numbers of the form

$$x = a \cdot 10^b$$

integer number with  $t$  digits

with  $0 \leq a < 1 \Rightarrow a = \pm 0. \overbrace{a_1 a_2 \dots a_t}$

and  $b = \pm \underbrace{b_1 b_2 \dots b_l}$

integer number with  $l$  digits

One calls  $a$  the "mantissa" and  $t$  the mantissa length

Example:  $30.40 = 0.3040 \cdot 10^2$   
 $= 0.0304 \cdot 10^3$

The representation is therefore not unique. An fp representation is called "normalized" if  $a_1 \neq 0$  (we only consider these in the following)

So: fp numbers represent a finite set  $A$  (since  $t, l < \infty$ ) of rational numbers

problem:  $A \neq \mathbb{R} \rightarrow$  chapter 1.2.2

### 1.2.2 Rounding errors

Main problem: If  $x, y \in A$  it is not guaranteed that

$$x \pm y, x \cdot y, x / y \in A$$

Example: ( $t=3, l=1$ )

$$0.300 \cdot 10^0 + 0.300 \cdot 10^{-4} = 0.3003 \notin A$$

Therefore, one defines for  $x \notin A$  the value  $rd(x) \in A$  the number that is "closest" to  $x$ , i.e.

$$|x - rd(x)| \leq |x - y| \quad \forall y \in A$$

satisfies. There are 2 cases:

i)  $x = \pm m \cdot 10^b \notin A$  since the mantissa has more than  $t$  digits.

Then  $rd(x)$  is the well-known rounding procedure, so

$$m = 0.m_1 \dots m_t m_{t+1} \dots$$

is

$$rd(x) = \pm m' \cdot 10^b$$

and

$$m' = \begin{cases} 0.m_1 \dots m_t & \text{if } 0 \leq m_{t+1} \leq 4 \\ 0.m_1 \dots (m_t + 1) & m_{t+1} > 4 \end{cases}$$

ii)  $x = \pm m \cdot 10^b \notin A$  since  $b$  has more than  $l$  digits, so e.g. for  $l=2$  is

$$\underbrace{(0.1 \cdot 10^{60})}_{\in A} \cdot \underbrace{(0.1 \cdot 10^{60})}_{\in A} = \underbrace{0.1 \cdot 10^{120}}_{\notin A}$$

This is called "exponent overflow" and is almost always the result of a programming error (not relevant here)

More important: since most operations result in numbers  $x \notin A$

there occur relative rounding errors of the order of magnitude

$$\left| \frac{\text{rd}(x) - x}{x} \right| \leq \epsilon$$

with the so-called "machine precision"  $\epsilon = 5 \cdot 10^{-(t+1)}$

### 1.2.3 fp arithmetic

Since rounding is required after many calculation operations, mathematically speaking, a new so-called fp arithmetic results in the set  $A$ :

Let  $x, y \in A$ , then

$$x \oplus y \equiv \text{rd}(x + y)$$

$$x \odot y \equiv \text{rd}(x \cdot y)$$

$$x \oslash y \equiv \text{rd}(x / y)$$

Problem: Almost all rules of real numbers are no longer valid, e.g. neither an associativity nor a distributive law.

Counterexample to the associative law: ( $t=2, l=1$ )

$$a = 0.52 \cdot 10^0$$

$$b = -0.51 \cdot 10^0$$

$$c = 0.50 \cdot 10^{-2}$$

$$\Rightarrow i) (a \oplus b) \oplus c = 0.10 \cdot 10^{-1} \oplus 0.50 \cdot 10^{-2} = 0.15 \cdot 10^{-1}$$

$$ii) a \oplus (b \oplus c) = 0.52 \cdot 10^0 \oplus (-0.51 \cdot 10^0) = 0.10 \cdot 10^{-1}$$

So: Be very careful when adding/subtracting numbers of different sizes and subtracting approximately equal numbers:

Example: Calculation of  $\int_0^{2\pi} dx \frac{\sin(x)}{x}$

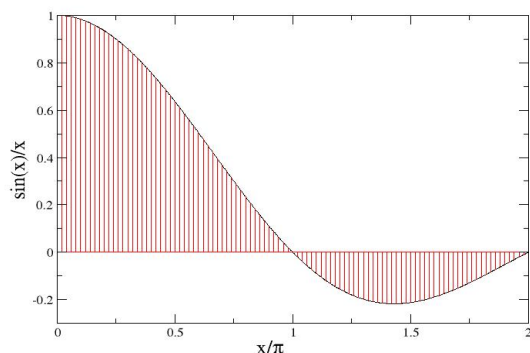
Simplest procedure ( $\rightarrow$  Chapter 4):

Break down the interval  $[0, 2\pi]$  into small intervals of width

Then

$$\Delta x = \frac{2\pi}{N}$$

$$\int_0^{2\pi} f(x) dx \approx \Delta x \sum_{i=1}^{n-1} f(i \cdot \Delta x)$$



Problem: If you naively just add up the strips, then with a large  $N$  you will eventually add up large numbers (the strips added up to that point) and small numbers (the next strip)  $\rightarrow$  big error.

### 1.3 Stability of iterative algorithms

Since there is always "numerical noise", algorithms must be stable with respect to this.

Example: Find numerical solution of the ordinary differential eq

$$\dot{\gamma} = -\gamma ; \gamma(0) = 1 \quad (1) \quad [\text{exact: } \gamma(t) = e^{-t}]$$

Anticipation of Chapter 7:

A) "Euler method" with discretization of  $t$ :

$$\tilde{\gamma}_n = \gamma(t_n) = \gamma(n \cdot \Delta t) \text{ with } \Delta t \ll 1$$

$$\Rightarrow \tilde{\gamma}_{n+1} - \tilde{\gamma}_n \stackrel{(\text{Taylor})}{=} \underbrace{\dot{\gamma}^{(1)}(t_n)}_{= -\tilde{\gamma}_n} \cdot \Delta t + O(\Delta t^2)$$

$$\Rightarrow \tilde{\gamma}_{n+1} = \tilde{\gamma}_n (1 - \Delta t) \quad (2)$$

so: ODE (1) is replaced by recursion equation (2)

Case 1: no rounding errors

⇒ solution of (2), with initial condition  $\gamma_0 = 1$  is  

$$\gamma_n = (1 - \Delta t)^n$$

⇒  $\gamma(t) = (1 - \frac{t}{n})^n \xrightarrow{n \rightarrow \infty} e^{-t}$

i.e. one gets an arbitrarily good approximation by increasing  $n$

Case 2: small rounding error at the beginning  $\gamma_0 = 1 + \Delta \gamma$   
 as above  
 ⇒  $\gamma(t) = (1 + \Delta \gamma) (1 - \frac{t}{n})^n \xrightarrow{n \rightarrow \infty} e^{-t} + \Delta \gamma e^{-t}$   
 $\lim_{t \rightarrow \infty} \frac{\Delta \gamma e^{-t}}{e^{-t}} = \Delta \gamma$

⇒ algorithm is stable

Remark: We have only shown stability for the simplest of all disturbances  $\gamma_0 \rightarrow \gamma_0 + \Delta \gamma$ . In reality there will be a  $\Delta \gamma_n$  in every step.

B) We now seemingly "improve" the Euler method by symmetrization

$$\gamma_{n+1} - \gamma_{n-1} \approx 2 \dot{\gamma}(t_n) \cdot \Delta t + \underbrace{O(\Delta t^3)}$$

(1)  
 ⇒  $\gamma_{n+1} = -2 \Delta t \gamma_n + \gamma_{n-1}$  (3) more precise than in A)

Solution Ansatz for (3):  $\gamma_n = (1 - T)^n$  ( $T$ : wanted)

Insert:  $(1 - T)^2 + 2 \cdot \Delta t (1 - T) - 1 = 0$

⇒  $(1 - T) = -\Delta t \pm \sqrt{\Delta t^2 + 1} \xrightarrow{\Delta t \ll 1} \approx -\Delta t \pm 1$

⇒ General solution of (3) for  $\Delta t \ll 1$ :

$$\gamma_n = 2 (1 - \Delta t)^n + \rho (-1 - \Delta t)^n \quad (4)$$

with arbitrary  $2, \rho$

⇒ We require two initial conditions here  $\gamma_0 = 1$

$\gamma_1 = (1 - \Delta t)$  (from A)

Insert (4) into initial condition:  $2 = 1, \rho = 0$

$\Rightarrow \tilde{y}_m = (1 - \Delta t)^m \xrightarrow{t \rightarrow \infty} \text{exact solution as in A)}$

But now again  $\tilde{y}_0 = 1 + \Delta t$ , such that  $\beta \neq 0$

Problem: The solution corresponding to  $\beta$

$$(-1 - \Delta t)^m = (-1)^m \underbrace{\left(1 + \frac{\Delta t}{m}\right)^m}_{\approx \infty}$$

diverges and oscillates, because

$$\ln(\lambda) = m \ln\left(\frac{\Delta t}{m}\right) \xrightarrow{m \rightarrow \infty} \infty \Rightarrow \text{Algorithm is unstable}$$

So: The accuracy of an algorithm usually says nothing about its stability



## 2) Solution of linear systems of equations, SVD

In CP, many algorithms require the solution of systems of linear equations (SLE) of the form

$$\tilde{A} \vec{x} = \vec{b}$$

with the  $M \times N$  matrix  $\tilde{A} = \begin{pmatrix} a_{11} & \dots & a_{1N} \\ \vdots & & \vdots \\ a_{M1} & \dots & a_{MN} \end{pmatrix}$

and the vectors  $\vec{x} = (x_1, \dots, x_N)^T$  (wanted)  
 $\vec{b} = (b_1, \dots, b_M)$  (given)

W.L.O.G.:  $a_{ij}, x_i, b_i \in \mathbb{R}$ . Complex LSE can be rewritten as real by splitting them into real and imaginary parts:

$$(\tilde{A}_R + i \tilde{A}_I) (\vec{x}_R + i \vec{x}_I) = \vec{b}_R + i \vec{b}_I$$

$$\Rightarrow \begin{aligned} \tilde{A}_R \vec{x}_R - \tilde{A}_I \vec{x}_I &= \vec{b}_R \\ \tilde{A}_I \vec{x}_R + \tilde{A}_R \vec{x}_I &= \vec{b}_I \end{aligned}$$

$$\Rightarrow \begin{pmatrix} \tilde{A}_R & -\tilde{A}_I \\ \tilde{A}_I & \tilde{A}_R \end{pmatrix} \begin{pmatrix} \vec{x}_R \\ \vec{x}_I \end{pmatrix} = \begin{pmatrix} \vec{b}_R \\ \vec{b}_I \end{pmatrix}$$

As is well known, there are (mathematically) three cases

I)  $\exists$  exactly one solution (only possible if  $M \geq N$ )

II)  $\exists$  infinitely many solutions ("underdetermined SLE")

III)  $\exists$  no solution ("overdetermined SLE")

$\Rightarrow$  In CP usually the result of a programming error

Note: In numerical calculations the three cases can become blurred, e.g.  $M=N$ , then the following applies:

$\tilde{A} \vec{x} = \vec{b}$  has a unique solution if and only if  $|\tilde{A}| \neq 0$

However: In numerical calculations there are no exact zeros



⇒ It is not always clear for which values of  $|\tilde{A}| \neq 0$  the numerical solution of the SLE is reasonable

One is on the safe side if

$|\tilde{A}| \sim$  typical order of magnitude of the matrix elements  $|a_{i,j}| \neq 0$

## 2.1 Case 1: LSE with unique solution

### 2.1.1 $M=N$

A) Note: The solution of an SLE requires in general  $2n^3$  operations. The various methods differ only in terms of  $2$ .

How better not to solve the problem?

i) Calculation of  $\tilde{A}^{-1}$  and  
 $\vec{x} = \tilde{A}^{-1} \vec{b}$

Exception: you need  $\tilde{A}^{-1}$  anyway (which is rarely the case)

ii) Gaussian method from HÖMa 1

B) The "LU decomposition" is usually very efficient and stable:

Assume:  $\tilde{A}$  is regular (i.e.  $|\tilde{A}| \neq 0$  and the LSE has a unique solution). Then you can write  $\tilde{A}$  as

$$\tilde{A} = \tilde{P} \cdot \tilde{L} \cdot \tilde{U}$$

where  $\tilde{L}(\tilde{U})$  is a lower (upper) triangular matrix, i.e.

$$\tilde{L} = \begin{pmatrix} l_{11} & 0 & 0 & \dots & 0 \\ l_{21} & l_{22} & & & \\ & & \ddots & & \\ & & & l_{nn} & \end{pmatrix}, \tilde{U} = \begin{pmatrix} u_{11} & u_{12} & & & \\ 0 & u_{22} & & & \\ & & \ddots & & \\ 0 & & & & u_{nn} \end{pmatrix}$$

and the orthogonal "pivoting matrix" ( $\tilde{P}^T = \tilde{P}^{-1}$ ) in each row/column contains exactly one 1.

Proof: lengthy and boring ( $\rightarrow$  Ref. 5)

Practical procedure  $\rightarrow$  maybe exercises

After the LU decomposition you can solve  $\tilde{A} \vec{x} = \vec{b}$  in 3 steps:

$$\tilde{P} (\tilde{L} (\tilde{U} \cdot \vec{x})) = \vec{b}$$

$$\Rightarrow \tilde{P} \vec{z} = \vec{b} \Rightarrow \vec{z} = \tilde{P}^T \cdot \vec{b} \quad (1)$$

$$\Rightarrow \tilde{L} \cdot \vec{y} = \vec{z} \quad (2) \quad \Rightarrow \quad \tilde{U} \cdot \vec{x} = \vec{y} \quad (3)$$

Advantage: Equations (2)/(3) are easy to solve, e.g. (2):

$$z_i = \sum_{j=1}^n l_{ij} y_j$$

$$\Rightarrow, i=1: z_1 = l_{11} y_1 \Rightarrow y_1 = z_1 / l_{11}$$

$$i=2: z_2 = l_{21} y_1 + l_{22} y_2 \Rightarrow y_2 = \frac{1}{l_{22}} (z_2 - l_{21} y_1)$$

$$\text{so: } y_i = \frac{1}{l_{ii}} \left( z_i - \sum_{j=1}^{i-1} l_{ij} y_j \right) \quad (4)$$

Equation (4) always makes sense, because  $l_{ii} \neq 0 \quad \forall i$

Reason: Note ii) below

There are several very efficient algorithms for determining the LU decomposition, implemented e.g. in "LAPACK" or "Eigen".

Remarks on the LU decomposition:

i) Advantage of the LU decomposition:

It only has to be carried out once ( $N^3$  effort) and can then be applied to any  $\vec{b}$ . (equation (4) only contains  $O(N^2)$  operations).

$$\text{ii) } |\tilde{A}| = \underbrace{|\tilde{P}|}_{= \pm 1} \cdot \underbrace{|\tilde{L}|}_{\prod_{i=2}^n l_{ii}} \cdot \underbrace{|\tilde{U}|}_{\prod_{i=1}^n u_{ii}} \quad (5)$$

$\Rightarrow$  would be  $l_{i,i} = 0$  or  $u_{i,i} = 0 \Rightarrow |\tilde{A}| = 0 \Rightarrow \tilde{A}$  not regular  $\hookrightarrow$

iii) Routines are common in which  $l_{i,i} = 1 \nmid$ ,  
(can always be set up)

iv) Inverse matrix

$$\tilde{A}^{-1} = \tilde{u}^{-1} \tilde{L}^{-1} \tilde{\rho}^T$$

$\tilde{u}^{-1}, \tilde{L}^{-1}$  are again easy to be determined, e.g.:  $\tilde{L}^{-1}$

$$\text{S.t. } (\tilde{L}^{-1})_{i,i} = \bar{l}_{i,i}$$

$$\tilde{L} \cdot \tilde{L}^{-1} = \tilde{I} \Rightarrow \sum_j l_{i,j} \bar{l}_{j,k} = \delta_{i,k}$$

$$\Rightarrow \sum_{j \leq i} l_{i,j} \bar{l}_{j,k} = \delta_{i,k}$$

$$\bullet i=1 : l_{1,1} \bar{l}_{1,k} = \delta_{1,k} \Rightarrow \bar{l}_{1,k} = \delta_{1,k} / l_{1,1}$$

$$\bullet i=2 : l_{2,1} \bar{l}_{1,k} + l_{2,2} \bar{l}_{2,k} \Rightarrow l_{2,k} = \frac{1}{l_{2,2}} (\delta_{2,k} - l_{2,1} \bar{l}_{1,k})$$

$$\bullet i=k+1 : \bar{l}_{i,k} = \frac{1}{l_{i,i}} (\delta_{i,k} - \sum_{j=1}^{i-1} l_{i,j} \bar{l}_{j,k}) \quad (O(n^2))$$

Remark:  $\bar{l}_{i,k} \neq 0$  only for  $k \leq i$  (proof: by induction)

$\Rightarrow \tilde{L}^{-1}$  is also a lower triangular matrix

$\tilde{u}^{-1}$  can be calculated analogously

Remarks:

For special matrices there may also be faster algorithms, examples:

i) symmetric, Hermitian, orthogonal, unitary matrices

ii) "band-diagonal matrices"

$$a_{i,j} = 0 \quad \forall i,j \quad |i-j| > m$$

$$\tilde{A} = \begin{pmatrix} \ddots & \ddots & \ddots & 0 \\ \vdots & \ddots & \ddots & \vdots \\ 0 & \ddots & \ddots & \vdots \end{pmatrix}$$

iii) "tridiagonal matrices" (will become important later)  
like (ii) with  $m=1$ , i.e.

$$\tilde{A} = \begin{pmatrix} a_{11} & a_{12} & 0 & 0 \\ a_{21} & a_{22} & a_{23} & \\ 0 & a_{32} & & \\ 0 & & & \ddots \\ & & & & a_{nn} \end{pmatrix}$$

iv) Sparse matrices

Matrices for which the number  $N'$  of  $a_{ij}$  with  $a_{ij} \neq 0$  is  
 $N' \ll N^2$

Remark: ii) and iii) are special cases of iv) if

$$N \gg 1 \text{ and } m \ll N$$

Otherwise one must analyze how the  $a_{ij} \neq 0$  are distributed (randomly, on the edge, in clusters, etc.)

$\Rightarrow$  The numerical solution of such SLE is a field in itself

2.1.2  $M > N$

So:  $\tilde{A} \cdot \vec{x} = \vec{b} \quad (1) :$

$$\begin{pmatrix} \vdots \\ \vdots \\ \vdots \end{pmatrix}_{\tilde{A}} \cdot \begin{pmatrix} \vdots \\ \vdots \\ \vdots \end{pmatrix}_{\vec{x}} = \begin{pmatrix} \vdots \\ \vdots \\ \vdots \end{pmatrix}_{\vec{b}}$$

In general there is no unique solution here.

Mathematical insertion:

Let  $\tilde{A}$  be an  $M \times N$  matrix. If the  $N \times N$  matrix

$$\tilde{\rho} \equiv \tilde{A}^T \cdot \tilde{A}$$

is invertible, then the matrix

$$\tilde{A}_p^{-1} \equiv \tilde{P}^{-1} \cdot \tilde{A}^T$$

is called the "Moore-Penrose-inverse" or "pseudo-inverse" of  $\tilde{A}$ , because

$$\tilde{A}_p^{-1} \cdot \tilde{A} = (\tilde{A}^T \cdot \tilde{A})^{-1} \cdot \tilde{A}^T \cdot \tilde{A} = \hat{I}$$

Remarks:

- i) If  $M=N$ , then obviously is  $\tilde{A}_p^{-1} = \tilde{A}^{-1}$
- ii)  $\tilde{A} \cdot \tilde{A}_p^{-1}$  is not defined for  $M \neq N$

Back to the problem:

$$\tilde{A} \cdot \vec{x} = \vec{b} \quad (1) \quad | \cdot \tilde{A}^T \quad \text{Apply the 1st + on both sides}$$

$$\Rightarrow \tilde{P} \cdot \vec{x} = \tilde{A}^T \cdot \vec{b} \equiv \vec{b}' \quad (2)$$

If (1) has a unique solution, then this is obviously also a solution of (2).

Remaining question: Is this solution also unique with respect to (2)?

Answer: yes

Because:  $\tilde{P}$  is regular, as can be seen as follows.

We define:

$$F(\vec{x}) \equiv \frac{1}{2} (\tilde{A} \vec{x} - \vec{b})^2 \quad (3)$$

This function has a global minimum at the solution  $\vec{x}_0$  of (1), where

$$F(\vec{x}_0) = 0$$

Evaluation of (3):

$$\begin{aligned} (\tilde{A} \vec{x} - \vec{b})^2 &= \left( \sum_i \tilde{A}_{i;}; x_i - b_i \right)^2 \\ &= \sum_{i,j} x_i \tilde{A}_{i;}; \tilde{A}_{j;}; x_j - 2 \sum_i b_i \tilde{A}_{i;}; x_i + b_i^2 \end{aligned}$$

Summation over i:

$$F(\vec{x}) = \underbrace{\vec{x}^T \tilde{A}^T \tilde{A} \vec{x}}_{=\tilde{p}} - 2 \underbrace{\vec{x}^T \tilde{A}^T \vec{b}}_{=\vec{b}'} + \vec{b}^2$$

$$= \vec{x}^T \tilde{p} \vec{x} - 2 \vec{x}^T \vec{b}' + \vec{b}^2 \quad (4)$$

If  $\tilde{p}$  was not regular, there would be a homogeneous solution  $\vec{x}_h$  of (2) for which

$$\tilde{p} \cdot \vec{x}_h = \vec{0}$$

$$\Rightarrow F(\vec{x}_0 + \lambda \vec{x}_h) \stackrel{(\tilde{p}^T = \tilde{p})}{=} 2\lambda \underbrace{\vec{x}_h^T (\tilde{p} \cdot \vec{x}_0 - \vec{b}')}_{=\vec{0}} + \lambda^2 \underbrace{\vec{x}_h^T \tilde{p} \vec{x}_h}_{=\vec{0}}$$

$$= 0 \quad \forall \lambda$$

This contradicts the assumption that  $\vec{x}_0$  is the global minimum of  $F$

Thus:

- If (1) has a unique solution  $\vec{x}_0$  then is  $\vec{x}_0$  the unique solution of (2).
- If we know (from whatever source, e.g. for physical reasons) that (1) has a unique solution, we can alternatively solve the smaller and quadratic problem (2) using the method from 2.1.1
- If (1) has no solution, then the solution (2) is the least bad in the sense of the function (3) ("least square minimum"), because with  $\vec{x}_0$  being a solution of (2)

$$\vec{0}_x = \frac{\partial_x F(\vec{x})}{\vec{x} = \vec{x}_0} = 2(\tilde{p} \vec{x}_0 - \vec{b}') = \vec{0}$$

Alternative: Chapter 2.2

## 2.2 Cases I)–III): The singular value decomposition (SVD)

With the SVD we have a method with which we

- i) can analyze which of the cases I)–III) applies
- ii) can determine solutions/solution-spaces/least-square-minima

### 2.2.1 On mathematics

Theorem: Every (!)  $M \times N$  matrix  $\tilde{A}$  can be written as

$$\tilde{A} = \tilde{U} \cdot \tilde{W} \cdot \tilde{V}^T \quad (1)$$

Here,

- i)  $\tilde{U}$  is an orthogonal  $M \times M$  matrix, i.e.

$$\tilde{U} = (\vec{u}_1, \dots, \vec{u}_M) = \begin{pmatrix} \vec{u}_1' \\ \vdots \\ \vec{u}_M' \end{pmatrix}, \quad \begin{aligned} \vec{u}_i \cdot \vec{u}_j &= \delta_{i,j} \\ \vec{u}_i' \cdot \vec{u}_j' &= \delta_{i,j} \end{aligned}$$

- ii)  $\tilde{V} = (\vec{v}_1, \dots, \vec{v}_N)$  and hence  $\tilde{V}^T$  is an orthogonal  $N \times N$ -Matrix

- iii)  $\tilde{W}$  is an  $M \times N$  matrix, for which

$$w_{i,j} = w_i \delta_{i,j}$$

where the "singular values"  $w_i \geq 0$  are the square-roots of the eigenvalues of  $\tilde{A}^T \cdot \tilde{A}$  ( $M \geq N$ ) or of  $\tilde{A} \cdot \tilde{A}^T$  ( $M \leq N$ )

$M > N$  :

$$\tilde{W} = \begin{pmatrix} w_1 & & & \\ & \ddots & & \\ & & w_N & \\ 0 & \dots & 0 & \\ \vdots & & \vdots & \\ 0 & \dots & 0 & \end{pmatrix}$$

$M < N$ :

$$\tilde{W} = \begin{pmatrix} w_1 & & 0 & \dots & 0 \\ & \ddots & \vdots & & \vdots \\ & & w_M & 0 & \dots & 0 \end{pmatrix}$$

Remarks:

- i) The right-hand side of (1) is unique, except for



- any (orthogonality-preserving) linear combinations of  $\{\vec{u}_i, \vec{u}_j\}$  or  $\{\vec{v}_i, \vec{v}_j\}$   $w_i = w_j$
- simultaneous permutation of  $\vec{u}_i, w_i, \vec{v}_i$

e.g.

$$\begin{aligned}\vec{u}_1 &\leftrightarrow \vec{u}_3 \text{ in } \tilde{u} \\ w_1 &\leftrightarrow w_3 \text{ in } \tilde{w} \\ \vec{v}_1 &\leftrightarrow \vec{v}_3 \text{ in } \tilde{v}\end{aligned}$$

ii) Algorithms for determining the SVD: Ref. 2 or 4.

We rely on LAPACK, Eigen, etc.

Examples:

i)  $M > N$

$$\tilde{A} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 0 \end{pmatrix} \Rightarrow \tilde{u} = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 & 0 & -1 \\ 0 & \sqrt{2} & 0 \\ 1 & 0 & 1 \end{pmatrix}, \tilde{w} = \begin{pmatrix} \sqrt{2} & 0 \\ 0 & 1 \\ 0 & 0 \end{pmatrix}, \tilde{v} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$

ii)  $M < N$

$$\hat{A} = \begin{pmatrix} 0 & 0 & 1 \\ 1 & 1 & 0 \end{pmatrix} \Rightarrow \tilde{u} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}, \tilde{w} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & \sqrt{2} & 0 \end{pmatrix}, \tilde{v} = \begin{pmatrix} 0 & \sqrt{2} & \sqrt{2} \\ 0 & \sqrt{2} & -\sqrt{2} \\ 1 & 0 & 0 \end{pmatrix}$$

iii)  $M = N$

$$\hat{A} = \begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix}, \tilde{u} = \begin{pmatrix} 0 & 1 \\ -1 & 0 \end{pmatrix}, \tilde{w} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}, \tilde{v} = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$$

So: Even for  $M = N$ , the SVD has nothing to do with the diagonalization of  $\tilde{A}$

Proof of the theorem (for simplicity, for  $M = N$ )

i)  $\hat{H} \equiv \tilde{A}^T \tilde{A}$  (or  $\tilde{A} \tilde{A}^T$ ) is symmetric  $\Rightarrow$  eigenvalues  $\lambda \in \mathbb{R}$

Let  $\vec{v}_i$  be an eigenvector of  $\hat{H}$  with the eigenvalue  $\lambda_i$ . Then

$$\begin{aligned}\vec{v}_i^T \cdot \tilde{A}^T \cdot \tilde{A} \cdot \vec{v}_i &= \lambda_i \underbrace{\vec{v}_i^T \cdot \vec{v}_i}_{\geq 0} \Rightarrow \lambda_i \geq 0 \Rightarrow w_i \equiv \sqrt{\lambda_i} \in \mathbb{R} \\ &\underbrace{= (\tilde{A} \vec{v}_i)^T}_{\geq 0} \quad \geq 0\end{aligned}$$

$\hat{A}$  has  $N$  orthogonal eigenvectors  $\hat{A} \vec{v}_i = \lambda_i \vec{v}_i$

and we define the matrices

$$\tilde{V} = (\vec{v}_1, \dots, \vec{v}_N)$$

and

$$\tilde{W} = \begin{pmatrix} w_1 & & \\ & \ddots & \\ & & w_N \end{pmatrix}$$

W.L.O.G.:  $\left. \begin{array}{l} w_i \neq 0 \quad \forall i \leq r \\ w_i = 0 \quad \forall i > r \end{array} \right\} \begin{array}{l} r=0 \text{ or } r=N \text{ is also} \\ \text{possible} \end{array}$

ii) What is  $\tilde{U}$ ?

For all  $i \leq r$  we define vectors  $\vec{u}_i \equiv \frac{\hat{A} \cdot \vec{v}_i}{w_i} \quad (2)$

These are  $\neq \vec{0}$  and orthonormal, because

$$\vec{u}_i^T \cdot \vec{u}_j = \frac{\vec{v}_i^T \hat{A}^T \hat{A} \vec{v}_j}{w_i \cdot w_j} = \frac{w_j^2}{w_i \cdot w_j} \quad \underbrace{\vec{v}_i^T \cdot \vec{v}_j}_{= \delta_{ij}} = \delta_{ij} \quad \checkmark$$

The vectors  $\vec{u}_1, \dots, \vec{u}_r$  can then be chosen arbitrarily such that  $\vec{u}_i \cdot \vec{u}_j = \delta_{ij} \quad \forall i, j$

$$\Rightarrow \tilde{U} \equiv (\vec{u}_1, \dots, \vec{u}_r)$$

iii) still to be shown:

$$\hat{A} = \tilde{U} \tilde{W} \tilde{U}^T \quad \text{or equivalently} \quad \tilde{W} = \tilde{U}^T \cdot \hat{A} \cdot \tilde{U}$$

$$\hat{A} \cdot \tilde{U} = (\hat{A} \vec{u}_1, \dots, \hat{A} \vec{u}_r)$$

where

$$\bullet \hat{A} \cdot \vec{u}_i \stackrel{(2)}{=} w_i \vec{u}_i \quad \forall i \leq r$$

$$\bullet \hat{A} \cdot \vec{u}_i = \vec{0} \quad \forall i > r$$

$$\text{because } (\hat{A} \cdot \vec{u}_i) \cdot (\hat{A} \cdot \vec{u}_j) = \vec{v}_i^T \hat{A}^T \hat{A} \vec{v}_j = w_j^2 \underbrace{\vec{v}_i^T \cdot \vec{v}_j}_{=0 \quad i > r} = 0$$

$$\Rightarrow \hat{A} \cdot \vec{u}_i = \vec{0} \quad \forall i > r$$

So

$$[\tilde{u}^T \tilde{A} \tilde{v}]_{i,j} = \left[ \begin{pmatrix} \vec{u}_1 \\ \vdots \\ \vec{u}_n \end{pmatrix} \cdot (w_1 \vec{u}_1, \dots, w_n \vec{u}_n) \right]_{i,j} = w_i \vec{u}_i \cdot \vec{u}_j$$

$$= w_i \delta_{i,j} = [\tilde{w}]_{i,j} \quad \checkmark$$

## 2.2.2 Reminder: Linear mappings

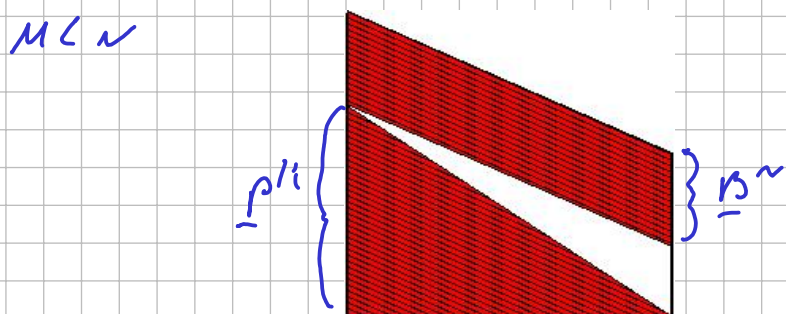
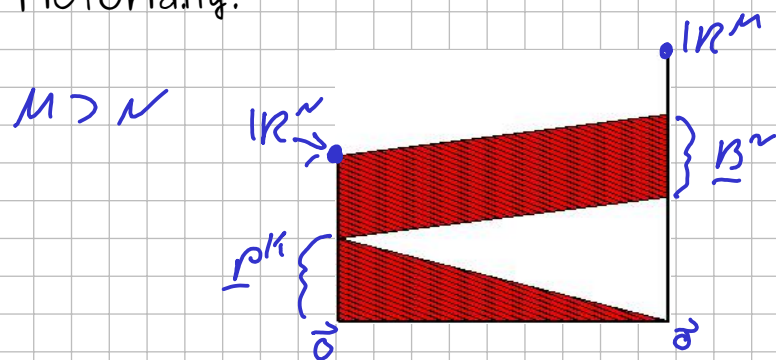
An  $M \times N$  matrix  $\tilde{A}$  defines a linear mapping

$$A: \mathbb{R}^N \rightarrow \mathbb{R}^M$$

A) Important statements and terms

- The image of  $A$  (or  $\tilde{A}$ ) is an  $r$ -dimensional subspace  $\underline{B}^r$  or  $\mathbb{R}^r$ .  $r$  is called the "rank" of the linear mapping (or of  $\tilde{A}$ ).
- The  $k$ -dimensional subspace  $\underline{P}^k$  of vectors  $\vec{x}$  in  $\mathbb{R}^N$  for which  $\tilde{A} \cdot \vec{x} = \vec{0}$  is called the "kernel" of the mapping. Its dimension  $k$  is called the "nullity".

iii) Pictorially:



iv) It applies  $r + k = N \quad (\Rightarrow r = N - k)$

Since  $r \leq m \Rightarrow N - k \leq m \Rightarrow k \geq N - m$

For  $M < N$ ,  $k$  is always  $> 0$

B) Application to LSE::

An LSE  $\tilde{A} \cdot \vec{x} = \vec{b}$  has

i) no solution if  $\vec{b} \notin \underline{B^r}$

ii) at least one solution if  $\vec{b} \in \underline{B^r}$ . This is unique if and only if  $k=0$ .

iii) If  $\vec{x}_0$  is a solution, then also  $\vec{x}_0 + \vec{x}_h$  with an arbitrary  $\vec{x}_h \in \underline{P^k}$  ( $\vec{x}_h$ : homogeneous solution)

### 2.2.3 Determination of $\underline{B^r}$ and $\underline{P^k}$ with the SVD

A) The columns  $\tilde{u}_i$  of  $\tilde{u}$ , with  $w_i \neq 0$  form an orthonormal basis (ONB) of  $\underline{B^r}$ , because:

We use the shorthand notation

$$\begin{aligned} \underline{B^r} &:= \tilde{A} \cdot \mathbb{R}^N = \text{span}(\tilde{A} \tilde{e}_1, \dots, \tilde{A} \tilde{e}_N) \\ &\quad \uparrow \\ &\quad \text{canonical basis} \\ &\equiv \text{space spanned by } \tilde{A} \tilde{e}_1, \dots, \tilde{A} \tilde{e}_N \end{aligned}$$

Then

$$\begin{aligned} \underline{B^r} &= \tilde{u} \tilde{w} \tilde{v}^T \cdot \mathbb{R}^N \\ &= \mathbb{R}^N \text{ (da } \tilde{v}^T \text{ orthogonal)} \end{aligned}$$

so the question remains: What is the image of  $\tilde{u} \cdot \tilde{w}$ ?

i)  $M \geq N$

$$\tilde{u} \cdot \tilde{w} = (\vec{u}_1, \dots, \vec{u}_M) \begin{pmatrix} w_1 & \dots & w_r & 0 & \dots & 0 \\ 0 & \dots & 0 & \dots & 0 & \dots & 0 \end{pmatrix} \begin{matrix} w_{r+1} \dots w_M \end{matrix}$$

$$(w_1 \vec{u}_1, \dots, w_r \vec{u}_r, \vec{0}, \dots, \vec{0}) \Rightarrow \text{claim}$$

ii)  $M < N$

$$\tilde{u} \cdot \tilde{w} = (\vec{u}_1 \dots \vec{u}_M) \begin{pmatrix} w_1 & & & 0 \\ & \ddots & & \\ & & w_r & \\ & & & \ddots \\ & & 0 & & 0 \end{pmatrix}$$

$$[w_1 \vec{u}_1, \dots, w_r \vec{u}_r, \vec{0}, \dots, \vec{0}] \Rightarrow \text{rank.}$$

B) The columns  $\vec{v}_i$  of  $\tilde{V}$  for which  $w_i \neq 0$  form an ONB of  $\underline{P}^k$  because:

i)  $M \leq N$

$$\begin{aligned} \tilde{w} \cdot \tilde{V}^T &= \tilde{w} (\vec{v}_1, \dots, \vec{v}_N)^T \\ &= (w_1 \vec{v}_1, \dots, w_M \vec{v}_M)^T \end{aligned}$$

[ $\vec{v}_{M+1}, \dots, \vec{v}_N$  have disappeared here]

so

$$(\tilde{w} \cdot \tilde{V}^T) \cdot \vec{v}_i = \begin{pmatrix} w_1 \vec{v}_1 \cdot \vec{v}_i \\ \vdots \\ w_i \vec{v}_i \cdot \vec{v}_i \\ \vdots \end{pmatrix} = \begin{pmatrix} 0 \\ \vdots \\ w_i \\ \vdots \end{pmatrix}$$

$$= w_i \vec{v}_i = \begin{cases} \vec{0} & \text{falls } w_i = 0 \Rightarrow \vec{v}_i \in \underline{P}^k \\ \neq 0 & \text{'' } w_i \neq 0 \Rightarrow \vec{v}_i \notin \underline{P}^k \end{cases}$$

ii)  $M > N$  : analogous

## 2.2.4 Analysis of linear systems of equations

A)  $M=N$ :

i) Case 1:  $w_i \neq 0 \quad \forall i$

$\Rightarrow \tilde{A} \vec{x} = \vec{b}$  has the unique solution

$$\vec{x}_0 = \tilde{A}^{-1} \cdot \vec{b} = \tilde{V} \cdot (\tilde{w}^{-1} \cdot (\tilde{u}^T \cdot \vec{b})) \quad (1)$$

with

$$\tilde{w}^{-1} = \begin{pmatrix} 1/w_1 & & 0 \\ & \ddots & \\ 0 & & 1/w_N \end{pmatrix}$$

Note: In (1) we are dealing with  $\sigma(n^2)$  operations. It therefore makes no sense to solve the SLE with another method after the SVD analysis.

ii) Whether  $\exists w_i = 0$  and  $\vec{b} \in \underline{B}^r \Rightarrow \exists \infty$  many solutions  
Whether  $\vec{b} \in \underline{B}^r$  can be easily checked: With  $\vec{u}_i (i=1, \dots, r)$  being a basis of  $\underline{B}^r$  one checks, if

$$\vec{b} - \sum_{i=1}^r (\vec{b} \cdot \vec{u}_i) \vec{u}_i = \vec{0}$$

(1) is also a solution here if  $\hat{w}^{-1}$  is set to  $\frac{1}{0} = 0$  namely the solution with minimum  $|\vec{x}|$ .

Proof: That it is a solution can be shown as in B)

to be shown:  $|\vec{x}_0 + \delta \vec{x}| \geq |\vec{x}_0|$

$\in \mathbb{R}^n$  d.h.  $\delta \vec{x} = \sum_i \lambda_i \vec{v}_i$ , where  $\lambda_i = 0$ , when  $w_i \neq 0$

$$\begin{aligned} & \stackrel{(\vec{v} \text{ orthogonal})}{=} |\vec{v} \hat{w}^{-1} \hat{u}^T \vec{b} + \delta \vec{x}| = |\vec{v} (\hat{w}^{-1} \hat{u}^T \vec{b} + \tilde{v}^T \cdot \delta \vec{x})| \\ & \stackrel{(\vec{v} \text{ orthogonal})}{=} |\underbrace{\hat{w}^{-1} \hat{u}^T \cdot \vec{b}}_{\vec{z}} + \underbrace{\tilde{v}^T \cdot \delta \vec{x}}_{\vec{h}}| \end{aligned}$$

$$\vec{h} = \begin{pmatrix} \vec{v}_1^T \\ \vdots \\ \vec{v}_n^T \end{pmatrix} \cdot \sum_i \lambda_i \vec{v}_i = \begin{pmatrix} \lambda_1 \\ \vdots \\ \lambda_n \end{pmatrix}$$

$$\text{so: } \left. \begin{array}{ll} \lambda_i \neq 0 \text{ only, when } w_i = 0 \\ \lambda_i \neq 0 & '' & w_i \neq 0 \end{array} \right\} \Rightarrow |\underbrace{\vec{z} + \vec{h}}_{=|\vec{x}_0 + \delta \vec{x}|}| \geq |\underbrace{\vec{z}}_{=|\vec{x}_0|}| \quad \checkmark$$

iii) Case 3:  $\exists w_i \neq 0$  und  $\vec{b} \notin \underline{B}^r \Rightarrow \exists$  no solution

In this case, (1) is the minimum of  $|\tilde{A} \cdot \vec{x} - \vec{b}|$

$\Rightarrow$  again a least-square-minimum (lsm) as in 2.1.2.

Proof: maybe exercises

B)  $M \neq N$

Analysis of the solution space and solution (or lsm) with (1) as in A)

Physically the most important case again: A seemingly over-determined LSE with a unique solution as in 2.1.2, i.e. here

- i)  $w_i \neq 0 \forall i \Rightarrow \underline{p}^N = \vec{0} \Rightarrow$  no hom. solution of the LSE
- ii)  $\vec{b} \in \mathbb{R}^r$

Then equation (1) applies again where

$$\tilde{W}^{-1} = \underbrace{\begin{pmatrix} 1/w_1 & & 0 & \dots & 0 \\ & \ddots & & & \\ & & 1/w_N & & 0 \\ & & & \ddots & \\ 0 & \dots & 0 & & 1/w_r \end{pmatrix}}_M \vec{v}$$

Proof: simply insert into SLE

$$\begin{aligned} \hat{A} \vec{x}_0 &= (\tilde{U} \cdot \tilde{W} \cdot \tilde{V}^T) (\tilde{V} \tilde{W}^{-1} \tilde{U}^T \vec{b}) \\ &= \tilde{U} \tilde{W} \tilde{W}^{-1} \tilde{U}^T \vec{b} \end{aligned}$$

It is

$$\tilde{W} \cdot \tilde{W}^{-1} = \begin{pmatrix} \tilde{\Gamma}_{N \times N} & 0 \\ 0 & 0 \end{pmatrix}$$

in A) a matrix  $\tilde{\Gamma}_{r \times r}$

If  $\vec{b} \in \mathbb{R}^r$  so is  $\vec{b} = \underbrace{\tilde{U}^T \vec{b}}_{\substack{N \text{ not } M \\ \text{not } M \times r}} \vec{u}_i$  and  $\tilde{U}^T \vec{b}$  has non-

vanishing components only in the first  $N$  rows. This means that  $\tilde{W} \cdot \tilde{W}^{-1}$  has the effect of an  $M \times M$   $\tilde{\Gamma}$  matrix

$$\Rightarrow \hat{A} \vec{x}_0 = \tilde{U} \tilde{U}^T \cdot \vec{b} = \vec{b} \quad \checkmark$$

## 2.2.5 Further applications of SVD

A) Construction of orthonormal bases

Given  $N$  vectors  $\vec{a}_i (i=1, \dots, N)$  in  $\mathbb{R}^M$

These form an  $N$ -dimensional subspace  $\mathbb{A}^N$  of  $\mathbb{R}^M$



Wanted: An ONB of  $\underline{A}'$

Textbook method: Gram-Schmidt method

→ Numerically very susceptible to rounding errors

much better: SVD of the matrix

$$\tilde{A} = (\vec{a}_1, \dots, \vec{a}_n) = \tilde{U} \cdot \tilde{W} \cdot \tilde{V}^T$$

Then the vectors  $\vec{u}_i$  with  $w_i \neq 0$  is the seeked ONB

Reason: The image of  $\tilde{A}$  is exactly the space  $\underline{A}'$  spanned by the  $\vec{a}_i$   $\xrightarrow{3.2.3}$  B.S.H.

B) Matrix approximation → data compression

Idea for  $M=N$  ( $M \neq N$  is analogous):

$$a_{ij} = \sum_{k=1}^N w_k u_{ik} v_{jk} \quad (w_{i+1} \leq w_i, w.l.o.z.)$$

$$\text{i.e. } \hat{A} = \sum_{k=1}^N w_k \vec{u}_k \vec{v}_k^T$$

Obviously the largest contributions come from the large  $w_k$

Approximation for  $\tilde{A}$ : Termination of the sum at  $r < N$

⇒ Only the  $\vec{u}_k, \vec{v}_k$  with  $k \leq r$  must be stored

$N=M=512$  pixels with  $a_{i,j}$  = grayscale

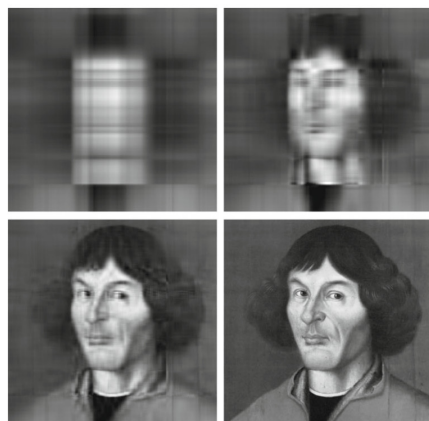


Fig. 3.7 Image compression by SVD. The original image is a  $512 \times 512$  matrix of pixels with 256 gray levels. Shown are the approximations of rank 2 (99.6% compression), rank 5 (99.9%), rank 15 (97.1%), and rank 100 (80.4%). See Problem 3.10.4

C) Face recognition (maybe exercise)

### 3) Eigenvalues (EW) and eigenvectors (EV)

#### 3.1 Mathematical repetition

Given:  $N \times N$  matrix  $\tilde{A}$

Wanted: All vectors  $\vec{v} \in \mathbb{R}^N$  with

$$\tilde{A} \cdot \vec{v} = \lambda \cdot \vec{v}, \quad \lambda \in \mathbb{C}$$

$$\Rightarrow (\tilde{A} - \lambda \tilde{I}) \cdot \vec{v} = \vec{0}$$

$$\Rightarrow P(\lambda) = |\tilde{A} - \lambda \tilde{I}| \stackrel{!}{=} 0$$

Mathematical textbook method:

i) Determine roots  $\lambda_i$  of  $P(\lambda)$

ii) solve for each  $\lambda_i$  the SLE

$$\tilde{A} \vec{v}_i = \lambda_i \vec{v}_i$$

However, this is numerically inefficient.

Alternative: Find a similarity transformation  $\tilde{Z}$ ,

$$\tilde{A} \rightarrow \tilde{A}' = \tilde{Z}^{-1} \cdot \tilde{A} \cdot \tilde{Z},$$

such that

$$\tilde{A}' = \begin{pmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_N \end{pmatrix}$$

is diagonal. Then the  $\lambda_i$  are the EV of  $\tilde{A}$  with eigenvectors

$$\vec{v}_i = \tilde{Z} \cdot \vec{l}_i \quad (= i\text{-th column of } \tilde{Z})$$

$\vec{l}_i$  (cartesian basis vectors)

because:  $\tilde{A} \cdot \vec{v}_i = \tilde{Z} \cdot \tilde{Z}^{-1} \cdot \tilde{A} \cdot \tilde{Z} \cdot \vec{l}_i = \lambda_i \cdot \vec{v}_i$

$$\begin{aligned} &= \tilde{A}' \cdot \vec{l}_i \\ &= \lambda_i \cdot \vec{l}_i \end{aligned}$$

We restrict ourselves to

i) real symmetric matrices  $\tilde{A}^T = \tilde{A}$

$\Rightarrow \lambda_i \in \mathbb{R}$ ,  $\tilde{Z}$  is orthogonal,  $\vec{v}_i \cdot \vec{v}_j = 0$ , if  $\lambda_i \neq \lambda_j$

ii) Hermitian matrices  $\tilde{A}^+ = \tilde{A}$  ↓  
 $\Rightarrow \lambda_i \in \mathbb{R}, \tilde{U} \text{ is unitary}, \vec{v}_i^+ \vec{v}_i = 0$

Here, too, we can restrict ourselves to case i), because

$$\tilde{A} \cdot \vec{v} = \lambda \cdot \vec{v} \Rightarrow (\tilde{A}_n + i \tilde{A}_p) (\vec{v}_n + i \vec{v}_p) = \lambda (\vec{v}_n + i \vec{v}_p)$$

$$\Rightarrow \begin{pmatrix} \tilde{A}_n & -\tilde{A}_p \\ \tilde{A}_p & \tilde{A}_n \end{pmatrix} \begin{pmatrix} \vec{v}_n \\ \vec{v}_p \end{pmatrix} = \lambda \begin{pmatrix} \vec{v}_n \\ \vec{v}_p \end{pmatrix}$$

Symmetric, it

$$\tilde{A}^+ = \tilde{A}_n^T - i \tilde{A}_p^T = \tilde{A}_n + i \tilde{A}_p = \tilde{A}$$

$$\Rightarrow -\tilde{A}_p = \tilde{A}_p^T \quad \checkmark$$

### 3.2 Jacobi rotation

Idea: Find the largest (absolute value) off-diagonal element  $a_{i,j}$  and determine  $\tilde{U}(i,j)$ , such that

$$\tilde{A}' = \tilde{U}(i,j)^T \tilde{A} \tilde{U}(i,j)$$

$$a'_{i,j} = 0 \text{ is satisfied}$$

Suitable approach:

$$\tilde{U}(i,j) = \begin{pmatrix} 1 & & & \\ & \ddots & & \\ & & \tilde{u}_{i,i} & \dots & \tilde{u}_{i,j} \\ & & \vdots & \ddots & \vdots \\ & & \tilde{u}_{j,i} & \dots & \tilde{u}_{j,j} \\ & & & \ddots & \\ & & & & 1 \end{pmatrix}; \begin{pmatrix} \tilde{u}_{i,i} & \tilde{u}_{i,j} \\ \tilde{u}_{j,i} & \tilde{u}_{j,j} \end{pmatrix} = \begin{pmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{pmatrix} \stackrel{=c}{=} \begin{pmatrix} c & s \\ -s & c \end{pmatrix}$$

$\Rightarrow \theta$  is to be determined:

$$\begin{pmatrix} c & -s \\ s & c \end{pmatrix} \begin{pmatrix} a_{i,i} & a_{i,j} \\ a_{j,i} & a_{j,j} \end{pmatrix} \begin{pmatrix} c & s \\ -s & c \end{pmatrix} \stackrel{=0}{=} \begin{pmatrix} a'_{i,i} & 0 \\ 0 & a'_{j,j} \end{pmatrix}$$

$$\Rightarrow (c^2 - s^2) a_{i,j} + s \cdot c (a_{i,i} - a_{j,j}) = 0$$

$$\Rightarrow \frac{c^2 - s^2}{sc} = \frac{a_{j,j} - a_{i,i}}{2 a_{i,j}} \stackrel{= \mathcal{N}}{=}$$

$$(t \equiv t_{\max}) \Rightarrow \pm \frac{\gamma}{2} \left( \frac{\gamma}{t} - t \right) \quad \text{quadratic eq.}$$

$t \text{ or } \frac{\gamma}{t}$

Solution:  $t = \frac{\gamma}{\lambda \pm \sqrt{1 + \lambda^2}} \Rightarrow c = \frac{1}{\sqrt{1 + t^2}}, s = c \cdot t$

To be clarified: Which sign?

Smallest intrusion, if  $|\theta| \leq \frac{\pi}{4} (\Rightarrow |t| \leq \frac{\gamma}{2})$

since then  $(a_{i,i} - a_{i,i}') + (a_{j,j} - a_{j,j}')$  is smaller than with  $\epsilon > \frac{\pi}{4}$

$$\Rightarrow \left. \begin{array}{l} \lambda > 0, \text{ "+" sign} \\ \lambda < 0, \text{ "-" sign} \end{array} \right\} \text{compact } t = \frac{\text{sign}(\lambda)}{|\lambda| + \sqrt{1 + \lambda^2}}$$

This results in the following algorithm:

- 1) Choose  $i, j$  such that  $|a_{i,j}|$  takes a maximum
- 2) Calculate  $c$  and  $s$  using the above formulas
- 3) Transformation  $\tilde{A}' = \tilde{\epsilon}(i,j)^T \cdot \tilde{A} \cdot \tilde{\epsilon}(i,j)$   
 $\tilde{\epsilon} \rightarrow \tilde{\epsilon} \cdot \tilde{\epsilon}(i,j)$  (start with  $\tilde{\epsilon} = \tilde{\eta}$ )
- 4) Check whether

$$\text{off}(\tilde{A}') \equiv \sum_{j \neq m} |a_{j,m}|^2$$

is below a certain limit, otherwise  $\tilde{A}' \rightarrow \tilde{A}$  and back to 1).

Remarks:

Since  $\text{off}(\tilde{A}') = \text{off}(\tilde{A}) - 2|a_{i,j}|^2 \quad (1)$

this value becomes smaller and smaller. Mathematically, this does of course not necessarily lead to convergence  $\rightarrow 0$ .

According to all experience, however, this is always the case (more on strict convergence proofs in Ref. 3)

Proof of (1):

a) It is  $\tilde{A}' = \tilde{A}$  except in the rows/columns  $i$  and  $j$ .

$$b) \sum_{l,m} a_{lm}^2 = \text{Tr}(\tilde{A}' \cdot \tilde{A}') = \text{Tr}(\tilde{A} \cdot \tilde{A}) = \sum_{l,m} a_{lm}^2$$

c) same argument as in b):

$$a_{ii}^2 + a_{jj}^2 + 2a_{ij}^2 = (a'_{ii})^2 + (a'_{jj})^2$$

Invariance of the trace  
under basis transformation.

$$\Rightarrow \text{off}(\tilde{A}') = \sum_{l,m} (a_{lm})^2 - \sum_i (a'_{ii})^2 \stackrel{b)}{=} \text{off}(\tilde{A}) + \underbrace{\sum_l a_{ll}^2 - a'_{ll}^2}_{= -2a_{ij}^2} \checkmark$$

ii) For large  $N$ , step 1 is time consuming. Then better fixed run through  $i, j$  in a certain order

iii) In general, convergence is slow  $\Rightarrow$  better 3.3

### 3.3 Diagonalization using the Housholder algorithm

The Housholder algorithm determines in  $N-2$  steps an orthogonal  $\tilde{Z}$ , with which  $\tilde{A}' = \tilde{Z}^T \tilde{A} \tilde{Z}$  has a tridiagonal form.

Step 1:

$$\tilde{Z}(1)^T \tilde{A} \tilde{Z}(1) = \begin{pmatrix} a_{11} & k_1 & 0 & \dots & 0 \\ k_1 & & & & \\ 0 & & \tilde{A}(1) & & \\ \vdots & & & & \\ 0 & & & & \end{pmatrix}$$

$$\text{with } \tilde{A}'(1) = \begin{pmatrix} a'_{22} & \dots & a'_{2N} \\ \vdots & & \\ a'_{N2} & \dots & a'_{NN} \end{pmatrix}$$

Step 2:  $\tilde{Z}(2)^T \tilde{Z}(1)^T \tilde{A} \tilde{Z}(1) \tilde{Z}(2)$

$$= \begin{pmatrix} a_{11} & k_1 & 0 & \dots & 0 \\ k_1 & a'_{22} & k_2 & 0 & \dots & 0 \\ 0 & & k_2 & & & \\ \vdots & & & \tilde{A}'(2) & & \\ 0 & & & & & \end{pmatrix} ; \tilde{A}'(2) = \begin{pmatrix} a'_{33} & \dots \\ \vdots & \end{pmatrix}$$

usw.  $\vdots$

Step  $N-2$ : tridiagonal form

Approach for  $\tilde{\mathbf{x}}(i)$ :

$$\tilde{\mathbf{x}}(i) = \begin{pmatrix} \tilde{x}_i \\ \vdots \\ 0 \\ \vdots \\ \tilde{\mathbf{S}}(i) \end{pmatrix} \quad (1)$$

with the  $(N-i) \times (N-i) = n \times n$  matrix

$$\tilde{\mathbf{S}}(i) \triangleq \hat{\mathbf{A}}_{n \times n} - 2 \vec{u}(i) \cdot \vec{u}(i)^T \quad (2)$$

$\vec{u}(i)$ :  $n$ -dimensional vector with  $|\vec{u}(i)| = 1$  (3)

Check orthogonality of  $\tilde{\mathbf{S}}(i)$ : (is it  $\tilde{\mathbf{S}}(i)^T = \tilde{\mathbf{S}}(i)$ )

$$\begin{aligned} \Rightarrow \tilde{\mathbf{S}}(i)^T \tilde{\mathbf{S}}(i) &= (\hat{\mathbf{A}} - 2 \vec{u} \vec{u}^T) (1 - 2 \vec{u} \vec{u}^T) \\ &= \hat{\mathbf{A}} - 4 \vec{u} \vec{u}^T + 4 \underbrace{\vec{u} \vec{u}^T \vec{u} \vec{u}^T}_{\stackrel{(3)}{=} \hat{\mathbf{A}}} = \hat{\mathbf{A}} \end{aligned}$$

Step  $i=1$ :

$$\tilde{\mathbf{x}}(1) = \begin{pmatrix} \tilde{x}_1 \\ \vdots \\ \vec{v} \\ \vdots \\ \tilde{\mathbf{S}}(1) \end{pmatrix}$$

$$\Rightarrow \tilde{\mathbf{x}}(1)^T \tilde{\mathbf{A}} \tilde{\mathbf{x}}(1) = \begin{pmatrix} a_{11} & (\tilde{\mathbf{S}}(1) \cdot \vec{v})^T & \dots & \dots \\ \vdots & \tilde{\mathbf{S}}(1) & \underbrace{\begin{pmatrix} a_{1n} & \dots \\ \vdots & a_{nn} \end{pmatrix}}_{\hat{\mathbf{A}}'(1)} & \tilde{\mathbf{S}}(1) \end{pmatrix}$$

with  $\vec{v} = \begin{pmatrix} a_{21} \\ a_{31} \\ \vdots \\ a_{n1} \end{pmatrix}$

Remaining question: How to choose  $\vec{u}(1)$  so that

$$(\tilde{\mathbf{S}}(1) \cdot \vec{v})^T = (h_1, 0, \dots, 0)$$

We write  $\vec{u}(1) = \frac{\vec{p}(1)}{|\vec{p}(1)|}$  (1) with arbitrary  $|\vec{p}(1)|$

approach:  $\vec{p}(1) = \vec{v} \pm |\vec{v}| \vec{e}_1$  (2)  $\Rightarrow |\vec{p}(1)|^2 = 2(|\vec{v}|^2 \pm |\vec{v}|)$  (3)

$$\Rightarrow \tilde{\mathbf{S}}(1) \cdot \vec{v} = \vec{v} - 2 \vec{u}(1) (\vec{u}(1) \cdot \vec{v}) = 2 \vec{p}(1) \cdot \vec{v} \quad (4)$$

$$v_1 = (v) \vec{v} = \frac{v^2 \pm |\vec{v}| v_1}{v^2 \pm |\vec{v}| v_1} (\vec{v} \pm |\vec{v}| \vec{e}_1) = \pm |\vec{v}| \vec{e}_1$$

therefore: Goal achieved

$$(\tilde{S}(1) \cdot \vec{v})^T = (v_1, 0, \dots, 0)$$

$$\text{mit } v_1 = \pm |\vec{v}|$$

Choose the sign as in  $-v_1 = -2v_1$  to avoid rounding errors in

$$v_1(1) = v_1 \pm |\vec{v}|$$

(for the conceivable case that  $v_1 \approx |\vec{v}|$ )

Steps 2 to  $N-2$  are completely analogous

⇒ tridiagonal matrix

$$\tilde{T} = \tilde{e}(N-1) \dots \tilde{e}(1) \tilde{A} \tilde{e}(1) \dots \tilde{e}(N-1)$$

Caution: When calculating  $\tilde{S}(1) \tilde{A} \tilde{S}(1)$  don't multiply the matrices (but for example)

$$(\vec{u} \cdot \vec{u}^T) \tilde{A} (\vec{u} \cdot \vec{u}^T) = \vec{u} \underbrace{[\vec{u}^T (\tilde{A} \vec{u})]}_{\substack{\text{vector} \\ \text{number}}} \vec{u}^T$$

After Tridiagonalising of  $\tilde{A}$ , one can diagonalise it with  $N-1$  further Jacobi rotations.

⇒ Goal is achieved

### 3.4 Matrix diagonalization (MD) in quantum mechanics

#### 3.4.1 Basic idea

In quantum mechanics, MD occurs when

- i) the Hilbert space is finite by definition, e.g. the "Heisenberg model"

$$\hat{H} = \sum_{i,j} T_{ij}; \quad \hat{S}_i \cdot \hat{S}_j; \quad ; \quad \hat{S}_i = \begin{pmatrix} \hat{S}_{i,x} \\ \hat{S}_{i,y} \\ \hat{S}_{i,z} \end{pmatrix}$$



on a finite chain with  $L$  sites  $(i,j)$ .

Then you can choose a basis  $|\varphi_i\rangle$  of the Hilbert space and diagonalise  $H_{e,x} \equiv \langle \varphi_i | \hat{H} | \varphi_j \rangle$ .

- ii)  $\hat{H}$  is approximately diagonalised in a properly chosen sub-space with a finite Basis  $|\varphi_i\rangle$ , e.g. a  $H_2$ -Molekül

$$\hat{H} = \sum_{i=1}^2 \hat{H}_{0,i} + \hat{V}_2$$

with  $\hat{H}_{0,i} = \frac{\hat{p}_i^2}{2m} - \frac{e^2}{|\vec{R}_i - \vec{r}_i|}$  ( $\vec{R}_i \in \mathbb{R}^3$ : Nuclei position fixed)

$$\hat{V}_2 = \frac{e^2}{|\vec{r}_1 - \vec{r}_2|} - \frac{e^2}{|\vec{R}_1 - \vec{r}_2|} - \frac{e^2}{|\vec{R}_2 - \vec{r}_1|}$$

eigenstates  $\varphi_{i,m}(\vec{r}_i)$  with eigenvalues  $E_m (m=1, \dots, \infty)$  of  $\hat{H}_{0,i}$  are known ( $\rightarrow$  Physics IV or Theoretical Physics II)

$\Rightarrow$  Form from the  $N$  low-energy states

$$\varphi_{1,m}(\vec{r}_1), \varphi_{2,m}(\vec{r}_2) (m=1, \dots, \nu)$$

two-electron states either in 1st quantization (Slater determinants) or in 2nd quantization (Fock states)

$\Rightarrow$  Diagonalize  $\hat{H}$  in this subspace

$\Rightarrow$  Approximation for the low-energy spectrum of  $\hat{H}$

### 3.4.2 Symmetries

of decisive importance is the skillful choice of the base  $|\varphi_i\rangle$

Example: Heisenberg model (with spin  $S=\frac{1}{2}$ ),  $L$  lattice sites  
basis of the Hilbert space:  $|\vec{S}\rangle \equiv |S_1, S_2, \dots, S_L\rangle$  with

$$S_i = \uparrow \text{ or } \downarrow$$

$\hat{H}$  commutes with  $\hat{S}^z = \sum_i \hat{S}_i^z$

$$\hat{S}^2 = \left( \sum_i \hat{S}_i \right)^2 \quad (\text{total spin})$$

The states  $|\vec{\sigma}\rangle$  are eigenstates of  $\hat{S}_z$

$$\hat{S}_z |\vec{\sigma}\rangle = \left( \sum_i \sigma_i \right) |\vec{\sigma}\rangle \equiv S_z(\vec{\sigma}) |\vec{\sigma}\rangle$$

$$\Rightarrow \langle \vec{\sigma} | [\hat{H}, \hat{S}_z] | \vec{\sigma}' \rangle = 0$$
$$= (S_z(\vec{\sigma}) - S_z(\vec{\sigma}')) \langle \vec{\sigma} | \hat{H} | \vec{\sigma}' \rangle$$

$$\Rightarrow \langle \vec{\sigma} | \hat{H} | \vec{\sigma}' \rangle = 0 \quad \text{if} \quad S_z(\vec{\sigma}) \neq S_z(\vec{\sigma}')$$

thus:  $\langle \vec{\sigma} | \hat{H} | \vec{\sigma}' \rangle$  is block-diagonal

$\Rightarrow$  Diagonalization is only required in the subspace with fixed  $S_z$

- dimension of the entire Hilbert space:  $D_H = 2^L$

- dimension of the subspace with fixed  $S_z$

Be  $N_\sigma$  the number of places with spin  $\sigma$

$$\Rightarrow \left. \begin{aligned} L &= N_\uparrow + N_\downarrow \\ S_z &= \frac{1}{2} (N_\uparrow - N_\downarrow) \end{aligned} \right\} \Rightarrow N_\uparrow = \frac{1}{2} (L + 2S_z)$$

$$\Rightarrow D_{\text{un}}^1(S_z) = \binom{L}{N_\uparrow} = \binom{L}{\frac{1}{2}(L + 2S_z)}$$

= number of possibilities to distribute the  $N_\uparrow \uparrow$  spins

e.g.:  $L=10 \Rightarrow D_H = 1024$   
 $D_{\text{un}}^1(10) = 252$

Question: In general, how can an optimal base (i.e. with maximum block diagonality of  $H_{i,j}$ ) be found?

$\Rightarrow$  Group theory (summer semester 2026)

### 3.5 The power method and the Lanczos algorithm

Hamiltonian matrices are often sparse (given a suitable choice of basis).

Example: Spin  $1/2$  Heisenberg model with nearest neighbor

coupling, we use ascent and descent operators

$$\hat{H} = J \sum_{\langle i,j \rangle} \left[ \hat{S}_i^x \hat{S}_j^x + \frac{1}{2} (\hat{S}_i^+ \hat{S}_j^- + \hat{S}_i^- \hat{S}_j^+) \right]$$

(i,j) = nearest neighbors

Base as above: States  $|\vec{\sigma}\rangle$

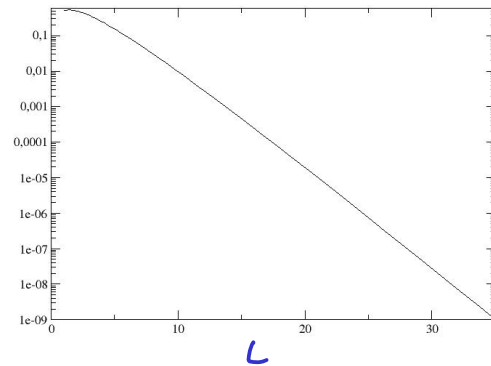
Each state  $|\vec{\sigma}\rangle$  has matrix elements  $\neq 0$  only with  $|\vec{\sigma}\rangle$ , and with  $|\vec{\sigma}'\rangle$  whose spins are flipped at exactly 2 neighboring sites

so: # of matrix elements  $\neq 0$  per row in  $\hat{H}$   $\approx L$

$\Rightarrow$  # of matrix elements  $\neq 0$  in  $\hat{H}$  :  $\approx L \cdot 2^L$

# of matrix elements:  $(2^L)^2$

$\Rightarrow$  Relation  $\frac{L \cdot 2^L}{(2^L)^2} = \frac{L}{2^L}$



To analyze the spectra of very large sparse matrices we consider two methods:

### 3.5.1 Power method

Algorithm for determining the largest eigenvalue (plus the associated eigenvector) of  $\tilde{A}$

- i) Choose a starting vector  $\vec{v}_0$
- ii) Iterate  $\vec{w}_m = \tilde{A} \vec{v}_{m-1}$  and  $\vec{v}_m = \frac{\vec{w}_m}{|\vec{w}_m|} = \frac{\tilde{A}^m \vec{v}_0}{|\tilde{A}^m \vec{v}_0|}$

converges to the searched eigenvector  $\vec{x}_\nu$  if

$\vec{v}_0 \cdot \vec{x}_\nu \neq 0$  and  $\lambda_\nu$  is not degenerate

Proof:

Let  $\vec{x}_i(\lambda_i)$  be the eigenvectors (eigenvalues) of  $\tilde{A}$

$$(\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_\nu)$$

Expand:

$$\vec{v}_0 = \sum_{i=1}^N \alpha_i \vec{x}_i$$

$$\Rightarrow \tilde{A}^m \vec{v}_0 = \sum_{i=1}^N \alpha_i \lambda_i^m \vec{x}_i = \lambda_N^m \left( \alpha_N \vec{x}_N + \sum_{i=1}^{N-1} \alpha_i \left( \frac{\lambda_i}{\lambda_N} \right)^m \vec{x}_i \right)$$

$$\Rightarrow \frac{\tilde{A}^m \vec{v}_0}{|\tilde{A}^m \vec{v}_0|} = \vec{x}_N + \underbrace{O\left(\left(\frac{\lambda_{N-1}}{\lambda_N}\right)^m\right)}_{\xrightarrow{m \rightarrow \infty} 0}$$

Remarks:

- i)  $\tilde{A} \Rightarrow -\tilde{A} \Rightarrow$  smallest eigenvalue of  $\tilde{A}$  with the same method
- ii) The numerical effort is low ( $O(N \cdot N^2)$ ) for sparse matrices still much smaller, but:

Convergence is often very slow, especially when

$$\left| \frac{\lambda_{N-1}}{\lambda_N} \right| \lesssim 1$$

- iii) With  $\lambda_N, \vec{x}_N$  found, one can define

$$\tilde{A}' = \tilde{A} - \lambda_N \vec{x}_N \cdot \vec{x}_N^T$$

$$\Rightarrow \tilde{A}' \cdot \vec{x}_N = 0 \cdot \vec{x}_N$$

With the same method one can determine  $\lambda_{N-1}, \vec{x}_{N-1}$  and in principle all  $\lambda_i, \vec{x}_i$

### 3.5.2 Krylov spaces

Wanted: Ground state  $\vec{v}_0$  of  $\tilde{A}$  ( $N \times N$ -matrix,  $\tilde{A}^T = \tilde{A}, \vec{v}_0 \in \mathbb{R}^N$ )

Principle way: Search for the minimum of

$$\mathcal{R}(\vec{v}) = \frac{\vec{v}^T \tilde{A} \vec{v}}{\vec{v}^T \cdot \vec{v}}, \quad \lambda_0 = \min[\mathcal{R}(\vec{v})] = \mathcal{R}(\vec{v}_0)$$

However, if performed exactly, this minimization would be numerically even more complex than the initial problem. The aim will be the minimization in subspaces

$$\underline{u}_k \equiv [\vec{p}_1, \dots, \vec{p}_k] \in \mathcal{S}_{\text{span}}(\vec{p}_1, \dots, \vec{p}_k)$$

of the

$$\vec{v}(\vec{d}) \sim \sum_{i=1}^k d_i \vec{p}_i \quad \vec{d} = \begin{pmatrix} d_1 \\ \vdots \\ d_k \end{pmatrix}$$

$$w_k = \min_{\vec{d}} \mathcal{J}_k(\vec{d}) = \min_{\vec{d}} \frac{\vec{v}(\vec{d})^T \hat{A} \vec{v}(\vec{d})}{\vec{v}(\vec{d})^T \vec{v}(\vec{d})} \Rightarrow d_{\min}$$

$$\Rightarrow \vec{v}^{(k)} \equiv \vec{v}(d_{\min}) = \text{optimum } \vec{v} \text{ in } \underline{u}_k$$

Obviously is

$$w_1 \geq w_2 \geq \dots \geq w_k \geq \lambda_0$$

Thus: By increasing  $k$ , one approaches the exact minimum of  $\mathcal{J}(\vec{v})$

Central question: How to choose  $\vec{p}_1, \dots, \vec{p}_k$  so that  $w_k$  is already a good approximation for  $\lambda_0$  even when  $k \ll N$

Inset:

anticipation of chapter 5: Minimization of functions  $\mathcal{J}(\vec{x})$

Simplest solution:

- 1) Choose a starting point  $\vec{x}_0$
- 2) Calculate  $\vec{F}(\vec{x}_0) = -\vec{\nabla} \mathcal{J} \big|_{\vec{x} = \vec{x}_0}$
- 3) Minimize  $\mathcal{J}(\vec{x}_0 - s \vec{F}(\vec{x}_0))$  (i.e. in the direction  $-\vec{F}(\vec{x}_0)$  with respect to  $s$ )
- 4) Back to 2)

Back to the initial question:

Let  $\vec{v}^k$  be the minimum of  $\mathcal{J}_k(\vec{d})$ . How do we then choose meaningfully  $\vec{p}_{k+1}$  (i.e.  $\underline{u}_{k+1}$ )

Answer: Such that  $\pm \vec{\nabla} \mathcal{J}(\vec{v}) \big|_{\vec{v} = \vec{v}^k} \in \mathcal{S}_{\text{span}}(\underline{u}_k, \vec{p}_{k+1})$ ,

because in this direction we should ideally change  $\vec{v}^k$

It is

$$\vec{0} \cdot \mathcal{L}(\vec{v}) \Big|_{\vec{v} = \vec{v}_k} \stackrel{\tilde{A}^T = \tilde{A}}{=} \frac{2}{(\vec{v}_k)^2} \left( \tilde{A} \cdot \vec{v}_k - \underbrace{\mathcal{L}(\vec{v}_k) \cdot \vec{v}_k}_{\in \mathcal{U}_k} \right)$$

contribution from the denominator.

$\Rightarrow$  Choose  $\vec{p}_{k+1} \sim \tilde{A} \cdot \vec{v}_k$  (1)

So:

- 1) start with (random)  $\vec{p}_1$
  - 2)  $\vec{p}_2 = \tilde{A} \cdot \vec{p}_1 \Rightarrow \vec{v}^{(2)} = \alpha_1 \vec{p}_1 + \alpha_2 \vec{p}_2$
  - 3) (1)  $\Rightarrow \vec{p}_3 = \tilde{A} \cdot \vec{v}^{(2)} = \alpha_1 \tilde{A} \vec{p}_1 + \alpha_2 \tilde{A} \vec{p}_2 = \alpha_1 \vec{p}_2 + \alpha_2 \tilde{A} \vec{p}_2$
- (according to the current idea by minimizing  $\mathcal{L}_2(\vec{v})$ )
- not required

but  $\alpha_1 \vec{p}_2$  is already in  $\mathcal{U}_2$

$\Rightarrow$  Instead of  $\vec{p}_3 = \tilde{A} \cdot \vec{v}^{(2)}$  we can also choose  $\vec{p}_3 = \tilde{A} \cdot \vec{p}_2$

$\Rightarrow$  The exact form of the optimal  $\vec{v}^{(2)}$  is therefore irrelevant for the iteration.

4) analogous for arbitrary  $k$   $\vec{p}_{k+1} = \tilde{A} \cdot \vec{p}_k$

The space

$$\mathcal{U}_k = \text{span}(\vec{p}_1, \vec{p}_2, \dots, \vec{p}_k) = \text{span}(\vec{p}_1, \tilde{A} \vec{p}_1, \dots, \tilde{A}^{k-1} \vec{p}_1)$$

is called "Krylov subspace".

### 3.5.3 The Lanczos method

Possible strategy:

Determine a Krylov space of sufficient size, construct the matrix  $H_{i,j} = \vec{p}_i^T \tilde{A} \cdot \vec{p}_j$  and determine its ground state

$\Rightarrow$  approximation for  $\lambda_0$

Problems:

- i) The vectors  $\vec{p}_i$  are generally not orthogonal
- ii) It is generally not even guaranteed that they are linearly independent  $\Rightarrow$  Dimension of  $\mathcal{U}_k$  possibly  $< k$

So: after the construction of a new  $\vec{p}_{k+1}$  it is better to orthonormalise ( $\vec{p}_{k+1} \rightarrow \vec{q}_{k+1}$ ) and abort, if  $\text{Dim}(\mathcal{H}_{k+1}) = k < k+1$

Problem: SVD does not make sense here, because

- i)  $\tilde{A}$  would have been in the base of the  $\vec{q}_i$ : few vanishing matrix elements (unlike the Lanczos method)
- ii) A full SVD would have to be calculated for each new vector  $\vec{p}_{k+1}$  i.e. without taking advantage of the fact that already  $\vec{q}_i \cdot \vec{q}_j = \delta_{ij}$   $\forall i, j \leq k$

The following algorithm is better:

(use bra-ket notation, i.e.  $\vec{q}_i \rightarrow |q_i\rangle$ ,  $\tilde{A} \rightarrow \hat{A}$   
 $\vec{q}_i^T \tilde{A} \vec{q}_i \rightarrow \langle \hat{A} \rangle_{q_i}$ )

starting values and vectors:  $|q_0\rangle = 0$

$|q_1\rangle = |q\rangle$  (arbitrary with  $\langle q|q\rangle = 1$ )

$$\gamma_1 = 1$$

$$\text{Iterations: } |q_{i+1}\rangle = \gamma_{i+1}^{-1} [(\hat{A} - \delta_i \hat{1}) |q_i\rangle - \gamma_i |q_{i-1}\rangle] \quad (1)$$

$$\equiv \gamma_{i+1}^{-1} |v_i\rangle$$

with  $\delta_i \equiv \langle \hat{A} \rangle_{q_i} \quad (2)$

$$\gamma_{i+1} \equiv \sqrt{\langle v_i | v_i \rangle} \quad (3)$$

Abort the iteration if  $\gamma_{i+1} = 0$  (i.e. when the dimension of the Krylov space no longer increases) or if the accuracy appears sufficient.

Proof of orthonormality:

i) Normalization: clear because of (3)

ii) Orthogonality by induction

A) Start of induction

i=2:  $\langle q_0 | q_2 \rangle = 0$  if  $|q_0\rangle = 0$  ✓



$$\gamma_2 |q_2\rangle = (\hat{A} - \delta_1 \hat{1}) |q_1\rangle = |v_2\rangle \quad (4)$$

$$\Rightarrow \langle q_1 | q_2 \rangle \sim (\langle A \rangle_{q_1} - \delta_1) \stackrel{(1)}{=} 0 \quad \checkmark$$

$$\text{and } \gamma_2^2 = \langle \hat{A}^2 \rangle_{q_1} - \langle \hat{A} \rangle_{q_1}^2 \quad (3)$$

$$i=3: \quad \gamma_3 |q_3\rangle = (\hat{A} - \delta_2 \hat{1}) |q_2\rangle - \gamma_2 |q_1\rangle$$

$$\Rightarrow \langle q_2 | q_3 \rangle \sim (\langle A \rangle_{q_2} - \delta_2) \stackrel{(2)}{=} 0$$

$$\langle q_1 | q_3 \rangle \sim (\underbrace{\langle q_2 | \hat{A} | q_2 \rangle}_{\stackrel{(4)}{=} \frac{1}{\gamma_2} \langle q_1 | \hat{A} (\hat{A} - \delta_1 \hat{1}) | q_1 \rangle} - \gamma_2) = 0$$

$$\stackrel{(4)}{=} \frac{1}{\gamma_2} \underbrace{\langle q_1 | \hat{A} (\hat{A} - \delta_1 \hat{1}) | q_1 \rangle}_{= \gamma_2^2}$$

(4)

B) Induction:

Let it now be  $\langle q_i | q_j \rangle = \delta_{ij}$   $\forall i, j \leq l$

too show:  $\langle q_i | q_{l+1} \rangle = 0 \quad \forall i \leq l$

a)  $\langle q_l | q_{l+1} \rangle = \langle q_{l-1} | q_{l+1} \rangle = 0$  as in (4)

b)  $\langle q_i | q_{l+1} \rangle$  for  $i < l-1$

$$\stackrel{(1)}{=} \gamma_{l+1}^{-1} \left( \underbrace{\langle q_i | \hat{A} | q_l \rangle}_{\approx 0} - \delta_l \underbrace{\langle q_i | q_l \rangle}_{=0} - \gamma_l \underbrace{\langle q_i | q_{l-1} \rangle}_{\approx 0} \right)$$

$$\hat{A}^T = \hat{A} \Rightarrow \langle q_l | \hat{A} | q_i \rangle \stackrel{(1)}{=} \gamma_i \underbrace{\langle q_l | q_i \rangle}_{=0} + \delta_i \underbrace{\langle q_l | q_i \rangle}_{=0} + \gamma_{i+1} \underbrace{\langle q_l | q_{i+1} \rangle}_{=0} = 0 \quad (5) \quad \checkmark$$

In the subspace  $\mathcal{U}_k$  and in the base of the  $|q_i\rangle$  is (due to (5))

$\hat{A}$  tridiagonal and has the matrix elements

$$\hat{T}_{i,i} \equiv \langle q_i | \hat{A} | q_i \rangle \stackrel{(1)}{=} \delta_i$$

$$\hat{T}_{(i+1),i} \equiv \langle q_{i+1} | \hat{A} | q_i \rangle = \gamma_i$$

$$\stackrel{1)}{=} \gamma_i \langle \eta_i | + \delta_i \langle \eta_{i-1} | + \gamma_{i+1} \langle \eta_{i+2} |$$

Therefore, the matrix to be diagonalized is

$$\tilde{T} = \begin{pmatrix} \delta_1 & \gamma_2 & 0 & \dots & 0 \\ \gamma_2 & \delta_2 & \gamma_3 & & \\ 0 & \gamma_3 & \delta_3 & & \\ & & & \ddots & \\ & & & & \end{pmatrix}$$

Summarized: (Lanczos method)

- 1) Determine the Lanczos basis  $\vec{\eta}_i$  for  $i \leq k$
- 2) Diagonalize  $\tilde{T}$ , as described in 3.2.
- 3) Convergence check, e.g.: is  $w_k \approx w_{k-1}$ ?

Remarks:

- i) Convergence is often very fast

Example: Heisenberg model,  $S = \hbar/2$ ,  $C = 0$

$\Rightarrow$  Hilbert space dimension  $D_H = 4^{10} \approx 10^7$

$$\Rightarrow \frac{w_{37} - w_{38}}{w_{38}} < 5 \cdot 10^{-8}$$

- ii) The orthogonalization formulas are identical to the Gram-Schmidt method, evaluated for our special vectors.

$$\vec{P}_i = \hat{A}^i \cdot \vec{\eta}$$

Problem with (1)-(3): rounding errors can compromise orthogonality for large  $k$ .

- iii) If  $k=N$  and  $\text{Sym}(\vec{\eta}_1, \dots, \vec{\eta}_N) = 1/N^N$

the Lanczos method is an alternative way to tridiagonalize a matrix.

But: The Housholder-Holder algorithm is more stable and

to be preferred.

iv) In the derivation of the Krylov spaces, we could also have replaced "min" by "max", with the same result.

→ At the same time, the Lanczos method also provides an approximation for the largest eigenvalue and its eigenvector.

v) The effort of the matrix-vector products  $\tilde{A} \cdot \vec{p}_k$  is  $O(N')$  where  $N'$  is the number of non-vanishing matrix elements. That means, that the dimension of the sparse matrix  $\tilde{A}$  is irrelevant and can therefore be very large.