

Numerical Methods for Computational Science and Engineering

Prof. R. Hiptmair, SAM, ETH Zurich

(with contributions from Prof. P. Arbenz and Dr. V. Gradinaru)

Autumn Term 2016

(C) Seminar für Angewandte Mathematik, ETH Zürich

URL: <http://www.sam.math.ethz.ch/~hiptmair/tmp/NumCSE/NumCSE16.pdf>

I. Computing with Matrices and Vectors

I.1. Fundamentals

I.1.1. Notation

$\mathbf{A} := \begin{bmatrix} a_{11} & \dots & a_{1m} \\ \vdots & \boxed{\quad} & \vdots \\ a_{n1} & \dots & a_{nm} \end{bmatrix}$

- entry $(\mathbf{A})_{i,j} = a_{ij}, 1 \leq i \leq n, 1 \leq j \leq m, a_{ij} \in \mathbb{K}$
- i -th row, $1 \leq i \leq n$: $a_{i,:} = (\mathbf{A})_{:,i}$
- j -th column, $1 \leq j \leq m$: $a_{:,j} = (\mathbf{A})_{:,j}$
- \mathbf{A} block → matrix block (sub-matrix) $(a_{ij})_{i=k,\dots,l}^{j=r,\dots,s} = (\mathbf{A})_{k:l,r:s}, 1 \leq k \leq l \leq n, 1 \leq r \leq s \leq m$

↑
column

I. Libraries

↳ Always use them for numerical linear algebra computations

I.2.3. Eigen

↳ Header-only C++ template library for numerical LA

Fundamental data type = matrix

Eigen::Matrix<typename Scalar, int nrows, int ncols>

fixed size matrices

dynamic matrices

C++11 code 1.2.12: Vector type and their use in EIGEN

```

1 #include <Eigen/Dense>
2
3 template<typename Scalar>
4 void eigenTypeDemo(unsigned int dim)
5 {
6     // General dynamic (variable size) matrices
7     using dynMat_t = Eigen::Matrix<Scalar, Eigen::Dynamic, Eigen::Dynamic>;
8     // Dynamic (variable size) column vectors
9     using dynColVec_t = Eigen::Matrix<Scalar, Eigen::Dynamic, 1>;
10    // Dynamic (variable size) row vectors
11    using dynRowVec_t = Eigen::Matrix<Scalar, 1, Eigen::Dynamic>;
12    using index_t = typename dynMat_t::Index;
13    using entry_t = typename dynMat_t::Scalar;
14
15    // Declare vectors of size 'dim'; not yet initialized
16    dynColVec_t colvec(dim); ↑ create vector objects
17    dynRowVec_t rowvec(dim);
18    // Initialisation through component access
19    for(index_t i=0; i< colvec.size(); ++i) colvec(i) = (Scalar)i;
20    for(index_t i=0; i< rowvec.size(); ++i) rowvec(i) = (Scalar)1/(i+1);
21    colvec[0] = (Scalar)3.14; rowvec[dim-1] = (Scalar)2.718;
22    // Form tensor product, a matrix, see Section 1.3.1
23    dynMat_t vecprod = colvec * rowvec;
24    const int nrows = vecprod.rows();
25    const int ncols = vecprod.cols();
26 }
```

[Matrix product]

② Component Access operator (int), (int,int),
→ Indices from 0!
Both Lvalue Rvalue

Line 23 :

tensor product

$$\begin{bmatrix} & \\ & \end{bmatrix} \cdot [] = \begin{bmatrix} & \\ & \end{bmatrix}$$

$$a, b \in \mathbb{R}^n : (\underline{a} \underline{b}^T)_{ij} = a_i b_j$$

Convenience data types : MatrixX*, * ∈ {i,f,cl,cd}
dynamic matrix

VectorX* ≈ dynamic column vector

C++11 code 1.2.14: Initializing special matrices in EIGEN

```

1 #include <Eigen/Dense>
2 // Just allocate space for matrix, no initialisation
3 Eigen::MatrixXd A(rows,cols);
4 // Zero matrix. Similar to matlab command zeros(rows,cols);
5 Eigen::MatrixXd B = MatrixXd::Zero(rows, cols);
6 // Ones matrix. Similar to matlab command ones(rows,cols);
7 Eigen::MatrixXd C = MatrixXd::Ones(rows, cols);
8 // Matrix with all entries same as value.
9 Eigen::MatrixXd D = MatrixXd::Constant(rows, cols, value);
10 // Random matrix, entries uniformly distributed in [0,1]
11 Eigen::MatrixXd E = MatrixXd::Random(rows, cols);
12 // (Generalized) identity matrix, 1 on main diagonal
13 Eigen::MatrixXd I = MatrixXd::Identity(rows,cols);
14 std::cout << "size of A = (" << A.rows() << ',' << A.cols() << ')' <<
    std::endl;
```

C++11 code 1.2.16: Demonstration code for access to matrix blocks in EIGEN

```

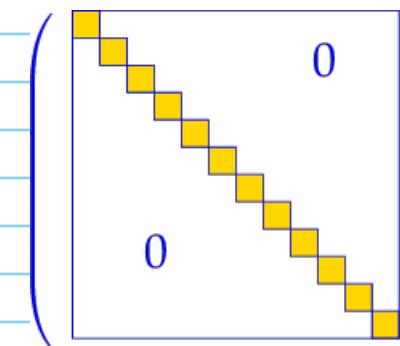
2 template<typename MatType>
3 void blockAccess(Eigen::MatrixBase<MatType> &M)
4 {
5     using index_t = typename Eigen::MatrixBase<MatType>::Index;
6     using entry_t = typename Eigen::MatrixBase<MatType>::Scalar;
7     const index_t nrows(M.rows()); // No. of rows
8     const index_t ncols(M.cols()); // No. of columns
9
10    cout << "Matrix M = " << endl << M << endl; // Print matrix
11    // Block size half the size of the matrix
12    index_t p = nrows/2, q = ncols/2;
13    // Output submatrix with left upper entry at position (i,j)
14    for(index_t i=0; i < min(p,q); i++)
15        cout << "Block (" << i << ',' << i << ',' << p << ',' << q
16        [read] << ") = " << M.block(i,i,p,q) << endl; ←
17    // l-value access: modify sub-matrix by adding a constant
18    M.block(1,1,p,q) += Eigen::MatrixBase<MatType>::Constant(p,q,1.0);
19    cout << "M = " << endl << M << endl;
20    // r-value access: extract sub-matrix
21    MatrixXd B = M.block(1,1,p,q);
22    cout << "Isolated modified block = " << endl << B << endl;
23    // Special sub-matrices
24    cout << p << " top rows of m = " << M.topRows(p) << endl;
25    cout << p << " bottom rows of m = " << M.bottomRows(p) << endl;
26    cout << q << " left cols of m = " << M.leftCols(q) << endl;
27    cout << q << " right cols of m = " << M.rightCols(p) << endl;
28    // r-value access to upper triangular part
29    const MatrixXd T = M.template triangularView<Upper>(); //
30    cout << "Upper triangular part = " << endl << T << endl;
31    // l-value access to upper triangular part
32    M.template triangularView<Lower>() *= -1.5; //
33    cout << "Matrix M = " << endl << M << endl;
34 }
```

Submatrix :
upper left corner
at (i,i) & size
 $p \times q$

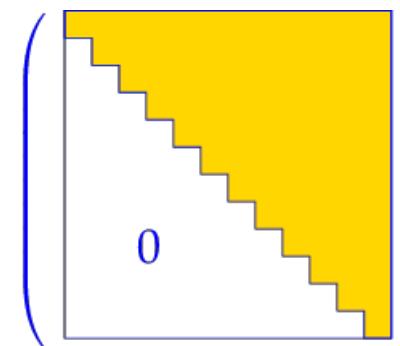
MATLAB triu

③

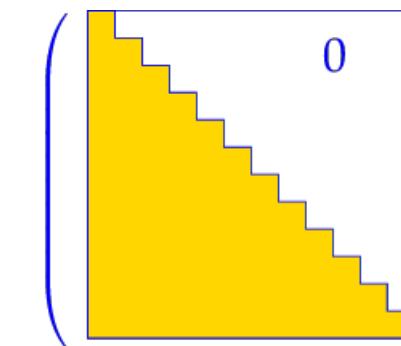
Triangular matrices :



diagonal matrix



upper triangular



lower triangular

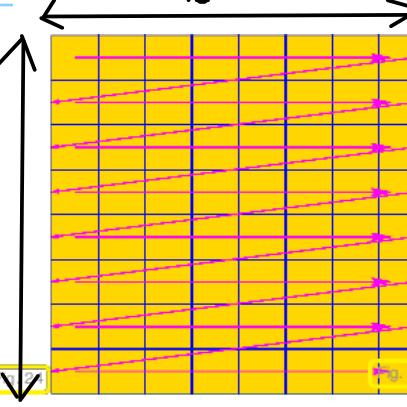
1.2.4. (Dense) Matrix Storage Formats

→ stored in one contiguous array

$$A \in \mathbb{K}^{m,n}$$

↔ Array of length $m \cdot n$

$$A = \begin{pmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{pmatrix}$$



Row major (C-arrays, bitmaps, Python):

A_arr	1	2	3	4	5	6	7	8	9
-------	---	---	---	---	---	---	---	---	---

Column major (Fortran, MATLAB, EIGEN):

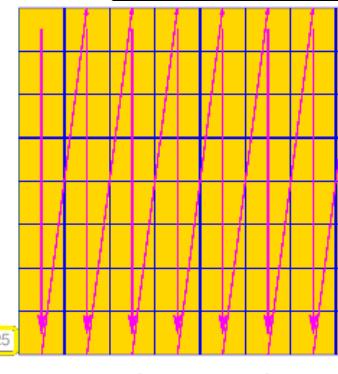
A_arr	1	4	7	2	5	8	3	6	9
-------	---	---	---	---	---	---	---	---	---

index mappings:

row major : [indexing from 0]

$$(i,j) \rightarrow i \cdot n + j$$

for $A \in \mathbb{R}^{m,n}$



C++11 code 1.2.22: Single index access of matrix entries in EIGEN

```

2 void storageOrder(int nrows=6,int ncols=7)
3 {
4     cout << "Different matrix storage layouts in Eigen" << endl;
5     // Template parameter ColMajor selects column major data layout
6     Matrix<double,Dynamic,Dynamic,ColMajor> mcm(nrows,ncols);
7     // Template parameter RowMajor selects row major data layout
8     Matrix<double,Dynamic,Dynamic,RowMajor> mmr(nrows,ncols);
9     // Direct initialization; lazy option: use int as index type
10    for (int l=1,i= 0; i< nrows; i++)
11        for (int j= 0; j< ncols; j++,l++)
12            mcm(i,j) = mmr(i,j) = l;
13
14    cout << "Matrix mmr = " << endl << mmr << endl;
15    cout << "mcm linear = ";
16    for (int l=0;l < mcm.size(); l++) cout << mcm(l) << ',';
17    cout << endl;
18
19    cout << "mmr linear = ";
20    for (int l=0;l < mmr.size(); l++) cout << mmr(l) << ',';
21    cout << endl;
22 }
```

Matrix-as-array access

Application : "reshaping" ↔ reinterpretation of matrix array data

Ex : Vectorization

$$A \in \mathbb{R}^{m,n} : \text{vec}(A) = \begin{bmatrix} (A)_{:,1} \\ (A)_{:,2} \\ \vdots \\ (A)_{:,n} \end{bmatrix}$$

C++11 code 1.2.28: Demonstration on how reshape a matrix in EIGEN

```

2 template<typename MatType>
3 void reshapeltest(MatType &M)
4 {
5     using index_t = typename MatType::Index;
6     using entry_t = typename MatType::Scalar;
7     const index_t nsize(M.size());
8
9     // reshaping possible only for matrices with non-prime dimensions
10    if ((nsize % 2) == 0) {
11        entry_t *Mdat = M.data(); // raw data array for M ← C-pointer
12        // Reinterpretation of data of M
13        Map<Eigen::Matrix<entry_t, Dynamic, Dynamic>> R(Mdat, 2, nsize / 2);
14        // (Deep) copy data of M into matrix of different size
15        Eigen::Matrix<entry_t, Dynamic, Dynamic> S =
16            Map<Eigen::Matrix<entry_t, Dynamic, Dynamic>>(Mdat, 2, nsize / 2);
17
18        cout << "Matrix M = " << endl << M << endl;
19        cout << "reshaped to " << R.rows() << 'x' << R.cols()
20        << " = " << endl << R << endl;
21        // Modifying R affects M, because they share the data space !
22        R *= -1.5;
23        cout << "Scaled (!) matrix M = " << endl << M << endl;
24        // Matrix S is not affected, because of deep copy
25        cout << "Matrix S = " << endl << S << endl;
26    }

```

```

1 Matrix M =
2 0 -1 -2 -3 -4 -5 -6
3 1 0 -1 -2 -3 -4 -5
4 2 1 0 -1 -2 -3 -4
5 3 2 1 0 -1 -2 -3
6 4 3 2 1 0 -1 -2
7 5 4 3 2 1 0 -1
8 reshaped to 2x21 =
9 0 2 4 -1 1 3 -2 0 2 -3 -1 1 -4 -2 0 -5 -3 -1 -6 -4 -2
10 1 3 5 0 2 4 -1 1 3 -2 0 2 -3 -1 1 -4 -2 0 -5 -3 -1

```

```

11 Scaled (!) matrix M =
12 -0 1.5 3 4.5 6 7.5 9
13 -1.5 -0 1.5 3 4.5 6 7.5
14 -3 -1.5 -0 1.5 3 4.5 6
15 -4.5 -3 -1.5 -0 1.5 3 4.5
16 -6 -4.5 -3 -1.5 -0 1.5 3
17 -7.5 -6 -4.5 -3 -1.5 -0 1.5

```

```

18 Matrix S =
19 0 2 4 -1 1 3 -2 0 2 -3 -1 1 -4 -2 0 -5 -3 -1 -6 -4 -2
20 1 3 5 0 2 4 -1 1 3 -2 0 2 -3 -1 1 -4 -2 0 -5 -3 -1

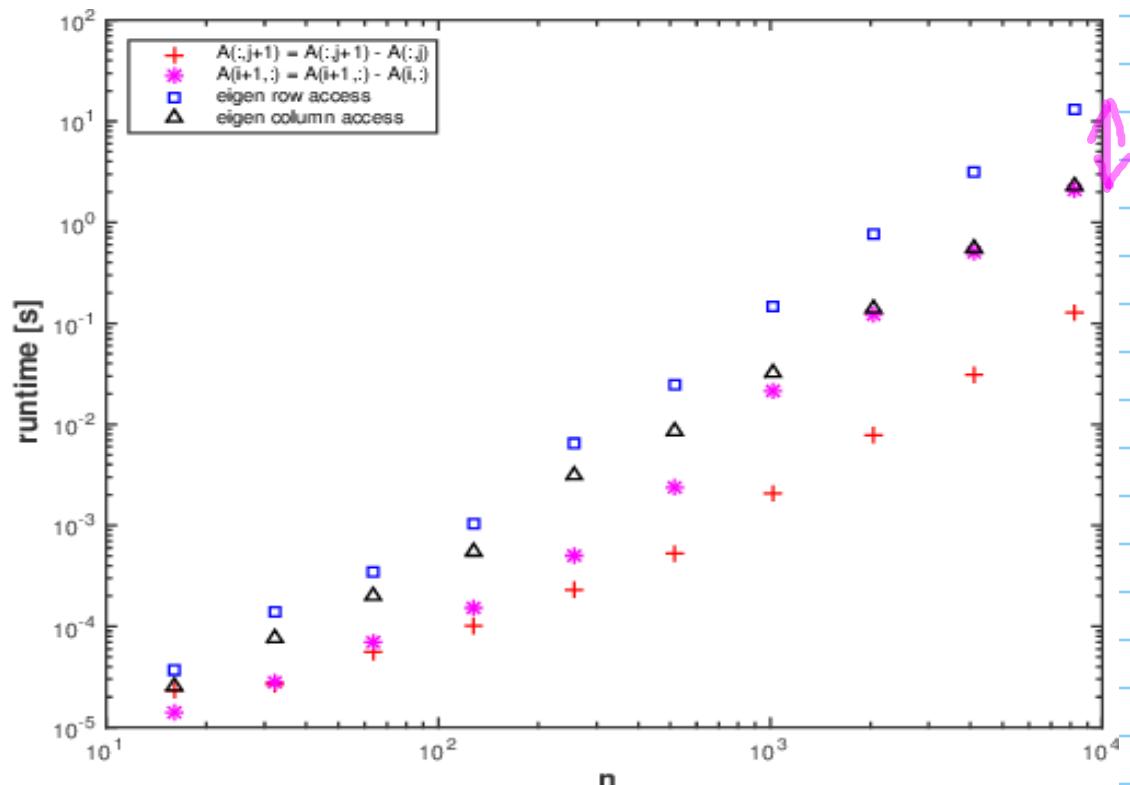
```

← C-pointer

Data Layout important for runtime performance !

Ex :

$$A_{\text{row}}(j+1) - = A_{\text{row}}(j)$$



F

Contiguous memory column access (for column Major storage layout) much faster : cache misses for row access.

5

I.4. Computational Effort

↳ No. of elementary operations '+', '−', '×', '÷'

"Computational effort $\not\propto$ runtime"



The computational effort involved in a run of a numerical code is only loosely related to overall execution time on modern computers.

Mainly determined by memory access patterns

I.4.1. Asymptotic Complexity

Definition 1.4.4. (Asymptotic) complexity

The **asymptotic complexity** of an algorithm characterises the **worst-case** dependence of its computational effort on one or more **problem size parameter(s)** when these tend to ∞ .

vector length, matrix dimensions

Notation : "Landau - O", e.g.

$$(*) \quad \text{Cost}(n) = O(n^\alpha), \alpha > 0$$

Strictly : $\Rightarrow \exists C > 0, n_0 \in \mathbb{N} : \text{Cost}(n) \leq Cn^\alpha \forall n \geq n_0$

tacit assumption of "sharpness":

(*) also means : $\text{Cost}(n) \neq O(n^\beta) \forall \beta < \alpha$

Asymptotic complexity $\not\propto$ runtime
 \Downarrow

dependence of runtime on problem size for large problems

$$\text{e.g. : } \text{cost}(n) = O(n^2)$$

$$\Rightarrow n \rightarrow n+2 : \text{cost} \times 4$$

Telling (polynomial) asymptotic complexity from runtime measurements

$$\text{conjecture : } t_i \propto C n_i^\alpha, i=1, \dots, N$$

$$\begin{matrix} \uparrow & \uparrow \\ \text{runtimes} & \text{problem sizes} \end{matrix}$$

$$\log t_i \approx \log C + \alpha \log n_i$$

Data points (n_i, t_i) roughly on a straight line with slope α in doubly logarithmic plot.

1.4.2. Cost of basic operations

operation	description	#mul/div	#add/sub	asympt. complexity
dot product	$(x \in \mathbb{R}^n, y \in \mathbb{R}^n) \mapsto x^H y$	n	$n - 1$	$O(n)$
tensor product	$(x \in \mathbb{R}^m, y \in \mathbb{R}^n) \mapsto xy^H$	nm	0	$O(mn)$
matrix product (*)	$(A \in \mathbb{R}^{m,n}, B \in \mathbb{R}^{n,k}) \mapsto AB$	mnk	$mk(n - 1)$	$O(mnk)$

↳ three loop implementation

1.4.3 Some tricks to reduce complexity

Ex: Exploit associativity

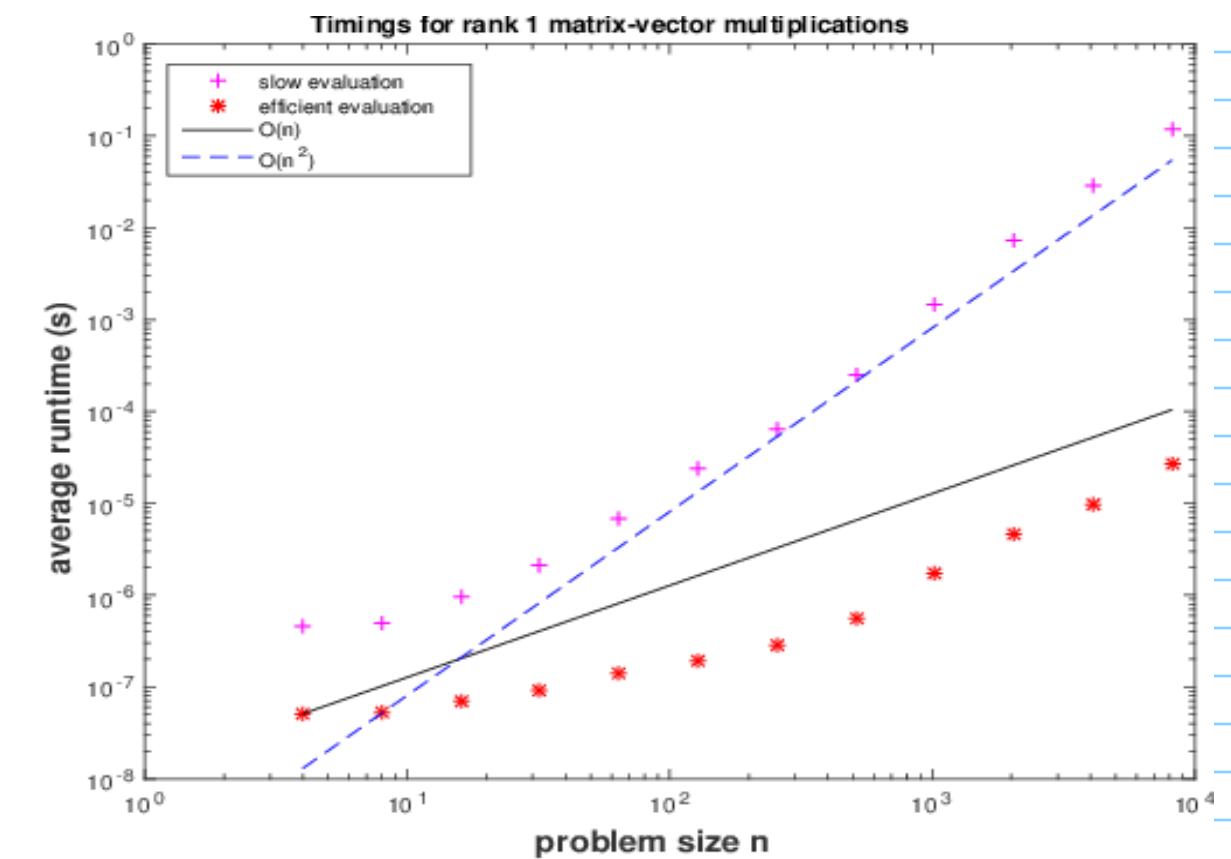
$$a \in \mathbb{R}^n, b \in \mathbb{R}^n, x \in \mathbb{R}^n$$

$$y = (ab^T)x. \quad (1.4.12)$$

`t = (a*b.transpose())*x;`

$$\begin{array}{c} \text{Diagram showing } t = a * b \cdot \text{dot}(x); \\ \text{Left side: } \begin{bmatrix} & & \\ & & \\ & & \end{bmatrix} \cdot \begin{bmatrix} & & \\ & & \\ & & \end{bmatrix} = \begin{bmatrix} & & \\ & & \\ & & n \times n \end{bmatrix} \cdot \begin{bmatrix} & & \\ & & \\ & & \end{bmatrix} \\ \text{Right side: } \begin{bmatrix} & & \\ & & \\ & & \end{bmatrix} \cdot \begin{bmatrix} & & \\ & & \\ & & \end{bmatrix} = \begin{bmatrix} & & \\ & & \\ & & \end{bmatrix} \cdot \begin{bmatrix} & & \\ & & \\ & & \end{bmatrix} \end{array}$$

(1.4.13):



Ex: "Hidden summation"

$$A \in \mathbb{R}^{n,p}, B \in \mathbb{R}^{n,p}, p \ll n$$

`Eigen::MatrixXd AB = A*B.transpose();
y = AB.triangularView<Eigen::Upper>()*x;`

$$Y = \text{triu}(A * B^T)$$

$\Rightarrow \text{Cost} = O(n^2)$ for $n \rightarrow \infty$, p fixed

(7)

 $p = 1 :$

$$\mathbf{y} = \text{triu}(\mathbf{ab}^T)\mathbf{x} = \begin{bmatrix} a_1 b_1 & a_1 b_2 & \dots & \dots & a_1 b_n \\ 0 & a_2 b_2 & a_2 b_3 & \dots & \dots & a_2 b_n \\ \vdots & \ddots & \ddots & \ddots & & \vdots \\ & & \ddots & \ddots & & \vdots \\ \vdots & & & \ddots & \ddots & \vdots \\ 0 & \dots & \dots & \dots & 0 & a_n b_n \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}$$

$$= \begin{bmatrix} a_1 & & & & \\ & \ddots & & & \\ & & \ddots & & \\ & & & \ddots & \\ & & & & a_n \end{bmatrix} \underbrace{\begin{bmatrix} 1 & 1 & \dots & \dots & 1 \\ 0 & 1 & 1 & \dots & 1 \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ & & \ddots & \ddots & \vdots \\ & & & \ddots & \vdots \\ \vdots & & & & \vdots \\ 0 & \dots & \dots & 0 & 1 \end{bmatrix}}_{\mathbf{T}} \begin{bmatrix} b_1 & & & & \\ & \ddots & & & \\ & & b_n & & \\ & & & \ddots & \\ & & & & x_n \end{bmatrix}$$

componentwise mult.

cumulative summation

componentwise mult.

Cost = $O(n)$ \Downarrow Cost = $O(n) [\cdot p]$

$$p > 1 : \quad \mathbf{AB}^T = \sum_{l=1}^p (\mathbf{A})_{:,e} (\mathbf{B})_{:,e}^T$$

$$\mathbf{y} = \text{triu}(\mathbf{AB}^T)\mathbf{x} = \sum_{l=1}^p \text{triu}((\mathbf{A})_{:,e} (\mathbf{B})_{:,e}^T)\mathbf{x}$$

Ex : Kronecker product

Definition 1.4.17. Kronecker product

The Kronecker product $\mathbf{A} \otimes \mathbf{B}$ of two matrices $\mathbf{A} \in \mathbb{K}^{m,n}$ and $\mathbf{B} \in \mathbb{K}^{l,k}$, $m, n, l, k \in \mathbb{N}$, is the $(ml) \times (nk)$ -matrix

$$\mathbf{A} \otimes \mathbf{B} := \begin{bmatrix} (\mathbf{A})_{1,1}\mathbf{B} & (\mathbf{A})_{1,2}\mathbf{B} & \dots & \dots & (\mathbf{A})_{1,n}\mathbf{B} \\ (\mathbf{A})_{2,1}\mathbf{B} & (\mathbf{A})_{2,2}\mathbf{B} & & & \vdots \\ \vdots & \vdots & & & \vdots \\ \vdots & \vdots & & & \vdots \\ (\mathbf{A})_{m,1}\mathbf{B} & (\mathbf{A})_{m,2}\mathbf{B} & \dots & \dots & (\mathbf{A})_{m,n}\mathbf{B} \end{bmatrix} \in \mathbb{K}^{ml,nk}.$$

Compute : $\mathbf{y} = (\mathbf{A} \otimes \mathbf{B})\mathbf{x}$ $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n,n}, \mathbf{x} \in \mathbb{R}^{n^2} \Rightarrow \text{cost}(n) = O(n^4)$

Smarter:

$$\mathbf{x} = \begin{bmatrix} \underline{x}^1 \\ \vdots \\ \underline{x}^n \end{bmatrix}, \quad \underline{x}^j \in \mathbb{R}^n$$

$$(\mathbf{A} \otimes \mathbf{B})\mathbf{x} = \begin{bmatrix} (\mathbf{A})_{1,1}\mathbf{Bx}^1 + (\mathbf{A})_{1,2}\mathbf{Bx}^2 + \dots + (\mathbf{A})_{1,n}\mathbf{Bx}^n \\ (\mathbf{A})_{2,1}\mathbf{Bx}^1 + (\mathbf{A})_{2,2}\mathbf{Bx}^2 + \dots + (\mathbf{A})_{2,n}\mathbf{Bx}^n \\ \vdots \\ \vdots \\ (\mathbf{A})_{m,1}\mathbf{Bx}^1 + (\mathbf{A})_{m,2}\mathbf{Bx}^2 + \dots + (\mathbf{A})_{m,n}\mathbf{Bx}^n \end{bmatrix}.$$

1.5. Machine Arithmetic

1.5.1. Loss of Orthogonalization

LA : Gram-Schmidt orthogonalization

Input: $\{a^1, \dots, a^k\} \subset \mathbb{K}^n$

```

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  $\ell = 1, 2, \dots, j-1$  do
5:      $q^j \leftarrow q^j - \langle a^j, q^\ell \rangle q^\ell$ 
6:   if ( $q^j = 0$ ) then STOP
7:   else {  $q^j \leftarrow \frac{q^j}{\|q^j\|_2}$  }
8: }
```

Output: $\{q^1, \dots, q^j\}$

If G.S. does not terminate prematurely *

- $(q^j)^T (q^i) = \delta_{ij}$

[orthonormal vectors]

- $\text{Span}\{q^1, \dots, q^j\}$
 $= \text{Span}\{a^1, \dots, a^j\}$

Proof \rightarrow LA

[* if $\{a^1, \dots, a^r\}$ linearly indep.]

Exp.: $A \in \mathbb{R}^{10,10}$, $(A)_{i,j} = \frac{1}{i+j-1}$, $1 \leq i, j \leq 10$
 "Hilbert matrix"

Test orthonormality of output vectors: $Q^T Q = I$

C++11 code 1.5.3: Gram-Schmidt orthogonalisation in EIGEN

```

2 template <class Matrix>
3 Matrix gramschmidt( const Matrix & A ) {
4   Matrix Q = A;
5   // First vector just gets normalized, Line 1 of (GS)
6   Q.col(0).normalize();
7   for(unsigned int j = 1; j < A.cols(); ++j) {
8     // Replace inner loop over each previous vector in Q with fast
9     // matrix-vector multiplication (Lines 4, 5 of (GS))
10    Q.col(j) -= Q.leftCols(j) * (Q.leftCols(j).transpose() *
11        A.col(j));//
12    // Normalize vector, if possible.
13    // Otherwise columns of A must have been linearly dependent
14    if( Q.col(j).norm() <= 10e-9 * A.col(j).norm() ) { //
15      std::cerr << "Gram-Schmidt failed: A has lin. dep
16      columns." << std::endl;
17      break;
18    } else { Q.col(j).normalize(); } // Line 7 of (GS)
19  }
20  return Q;
}
```

C++11 code 1.5.7: Wrong result from Gram-Schmidt orthogonalisation EIGEN

```

2 void gsroundoff(MatrixXd& A){
3   // Gram-Schmidt orthogonalization of columns of A, see Code 1.5.3
4   MatrixXd Q = gramschmidt(A);
5   // Test orthonormality of columns of Q, which should be an
6   // orthogonal matrix according to theory
7   cout << setprecision(4) << fixed << "I = "
8   << endl << Q.transpose()*Q << endl;
9   // EIGEN's stable internal Gram-Schmidt orthogonalization by
10  // QR-decomposition, see Rem. 1.5.9 below
11  HouseholderQR<MatrixXd> qr(A.rows(), A.cols()); //
12  qr.compute(A); MatrixXd Q1 = qr.householderQ(); //
13  // Test orthonormality
14  cout << "I1 = " << endl << Q1.transpose()*Q1 << endl;
15  // Check orthonormality and span property (1.5.2)
16  MatrixXd R1 = qr.matrixQR().triangularView<Upper>();
17  cout << scientific << "A-Q1*R1 = " << endl << A-Q1*R1 << endl;
18 }
```

Eigen's G.S.

II =

```

1.0000 0.0000 -0.0000 0.0000 -0.0000 0.0000 -0.0000 -0.0000 -0.0000 -0.0000
0.0000 1.0000 -0.0000 0.0000 -0.0000 0.0000 -0.0000 -0.0000 -0.0000 -0.0000
-0.0000 -0.0000 1.0000 0.0000 -0.0000 0.0000 -0.0000 -0.0000 -0.0000 -0.0000
0.0000 0.0000 0.0000 1.0000 -0.0000 0.0000 -0.0000 -0.0000 -0.0000 -0.0000
-0.0000 -0.0000 -0.0000 -0.0000 1.0000 0.0000 -0.0008 -0.0007 -0.0007 -0.0006
0.0000 0.0000 0.0000 0.0000 1.0000 -0.0540 -0.0430 -0.0360 -0.0289 -0.0289
-0.0000 -0.0000 -0.0000 -0.0008 -0.0540 1.0000 0.9999 0.9998 0.9996 -0.0000
-0.0000 -0.0000 -0.0000 -0.0007 -0.0430 0.9999 1.0000 1.0000 0.9999 -0.0000
-0.0000 -0.0000 -0.0000 -0.0007 -0.0360 0.9998 1.0000 1.0000 1.0000 -0.0000
-0.0000 -0.0000 -0.0000 -0.0006 -0.0289 0.9996 0.9999 1.0000 1.0000 -0.0000

```

No orthonormality ↑

► Computers cannot compute properly (in \mathbb{R})!

↳ Inevitably so, because they can handle only
finitely many numbers

-> Machine number (set M)

$M \subset \mathbb{R}$: finite & discrete

\exists largest/smallest (in modulus) machine numbers

Defined in Standard IEEE754/
IEC559 [for double/float]

underflow \ overflow can happen

C++11 code 1.5.21: Querying characteristics of double numbers

```

2 #include <limits>
3 #include <iostream>
4 #include <iomanip>
5
6 using namespace std;
7
8 int main() {
9     cout << std::numeric_limits<double>::is_iec559 << endl
10    << std::defaultfloat << numeric_limits<double>::min() << endl
11    << std::hexfloat << numeric_limits<double>::min() << endl
12    << std::defaultfloat << numeric_limits<double>::max() << endl
13    << std::hexfloat << numeric_limits<double>::max() << endl;
14 }

```

Output:

- 1 true
- 2 2.22507e-308
- 3 0010000000000000
- 4 1.79769e+308
- 5 7fffffffffffffff

$op : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ $op \in \{ +, -, \times, / \}$

[M is not closed under elementary arithmetic operations]

On computer $\tilde{op} : M \times M \rightarrow M$

→ implementation : $\tilde{op} = rd \circ op$

Definition 1.5.27. Correct rounding [to nearest machine number]

Correct rounding ("rounding up") is given by the function

$$rd : \begin{cases} \mathbb{R} \rightarrow M \\ x \mapsto \max_{\tilde{x} \in M} \arg \min_{\tilde{x} \in M} |x - \tilde{x}| \end{cases}$$

(12)

C++11 code 1.5.23: Demonstration of roundoff errors

```

2 #include <iostream>
3 int main(){
4     std::cout.precision(15);
5     double a = 4.0/3.0, b = a-1, c = 3*b, e = 1-c;
6     std::cout << e << std::endl;
7     a = 1012.0/113.0; b = a-9; c = 113*b; e = 5+c;
8     std::cout << e << std::endl;
9     a = 83810206.0/6789.0; b = a-12345; c = 6789*b; e = c-1;
10    std::cout << e << std::endl;
11 }

```

Output:

- 1 2.22044604925031e-16
- 2 6.75015598972095e-14
- 3 -1.60798663273454e-09

Controlling roundoff error :

Maximal relative error of rounding :

$$\text{EPS} = \max_{x \in \mathbb{R}} \frac{|rd(x) - x|}{|x|}$$

[overflow/underflow ignored]

[rel. error of \tilde{x} w.r.t. x = $\frac{|\tilde{x}-x|}{|x|} = \delta : \tilde{x} = x(1+\epsilon)$
 $|\epsilon| < \delta$.]

Assumption 1.5.32. "Axiom" of roundoff analysis

There is a small positive number **EPS**, the **machine precision**, such that for the elementary arithmetic operations $\star \in \{+, -, \cdot, /\}$ and "hard-wired" functions* $f \in \{\exp, \sin, \cos, \log, \dots\}$ holds

$$x \tilde{\star} y = (x \star y)(1 + \delta), \quad \tilde{f}(x) = f(x)(1 + \delta) \quad \forall x, y \in \mathbb{M},$$

with $|\delta| < \text{EPS}$.

↑ worst case estimate

Exact: $e = 0$

Computing in M

C++11-code 1.5.34: Finding out EPS in C++

```

2 #include <iostream>
3 #include <limits> // get various properties of arithmetic types
4 int main() {
5     std::cout.precision(15);
6     std::cout << std::numeric_limits<double>::epsilon() << std::endl;
7 }

```

Output:

1 2.22044604925031e-16

```

cout.precision(25);
double eps =
    numeric_limits<double>::epsilon();
cout << fixed << 1.0 + 0.5*eps << endl
    << 1.0 - 0.5*eps << endl
    << (1.0 + 2/eps) - 2/eps << endl;

```

Output:

1 1.0000000000000000000000000000000
2 0.999999999999998889776975
3 0.0000000000000000000000000000000

$\overbrace{1 + \delta}^{\approx 1} = 1$ compatible with "Axiom" ?
 $\Rightarrow \delta = (1 + \delta)(1 + \delta) - 1 = \delta^2$

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

$$\Rightarrow |\delta| = \left| \frac{1}{1 + \delta} \right| \leq \frac{1}{1 + \delta} = \frac{1}{2} \text{EPS}$$

OK ✓

13

Remark: Testing == 0.0

C++11 code 1.5.3: Gram-Schmidt orthogonalisation in EIGEN

```

2 template <class Matrix>
3 Matrix gramschmidt( const Matrix & A ) {
4     Matrix Q = A;
5     // First vector just gets normalized, Line 1 of (GS)
6     Q.col(0).normalize();
7     for(unsigned int j = 1; j < A.cols(); ++j) {
8         // Replace inner loop over each previous vector in Q with fast
9         // matrix-vector multiplication (Lines 4, 5 of (GS))
10        Q.col(j) -= Q.leftCols(j) * (Q.leftCols(j).transpose() *
11            A.col(j));//
12        // Normalize vector, if possible.
13        if( Q.col(j).norm() <= 10e-9 * A.col(j).norm() ) { //
14            std::cerr << "Gram-Schmidt failed: A has lin. dep
15            columns." << std::endl;
16            break;
17        } else { Q.col(j).normalize(); } // Line 7 of (GS)
18    }
19    return Q;
}

```

a numerical crime for results of floating point computations

replace with test of "relative smallness"

?

1.5.4. Cancellation

Ex: Computing the roots of a quadratic polynomial

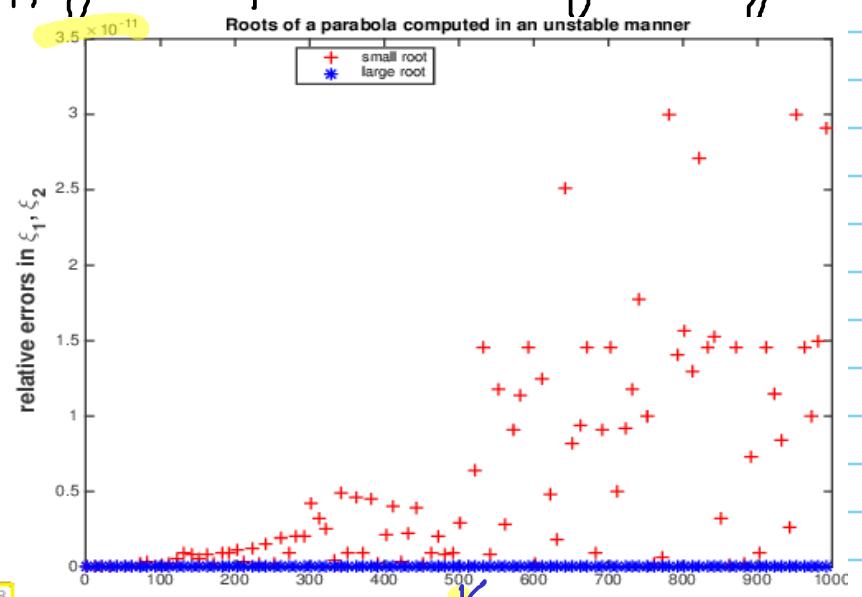
C++11-code 1.5.41: Discriminant formula for the real roots of $p(\xi) = \xi^2 + \alpha\xi + \beta$

```

2     ///! C++ function computing the zeros of a quadratic polynomial
3     ///!  $\xi \rightarrow \xi^2 + \alpha\xi + \beta$  by means of the familiar discriminant
4     ///! formula  $\xi_{1,2} = \frac{1}{2}(-\alpha \pm \sqrt{\alpha^2 - 4\beta})$ . However
5     ///! this implementation is vulnerable to round-off! The zeros are
6     ///! returned in a column vector
7     Vector2d zerosquadpol(double alpha, double beta){
8         Vector2d z;
9         double D = std::pow(alpha,2) - 4*beta; // discriminant
10        if(D < 0) throw "no real zeros";
11        else{
12            // The famous discriminant formula
13            double wD = std::sqrt(D);
14            z << (-alpha-wD)/2, (-alpha+wD)/2; //
15        }
16    }
17
}

```

$$\text{Apply to } p(\xi) = (\xi-\gamma)(\xi-\frac{1}{\gamma}) = \xi^2 - (\gamma + \frac{1}{\gamma})\xi + 1$$

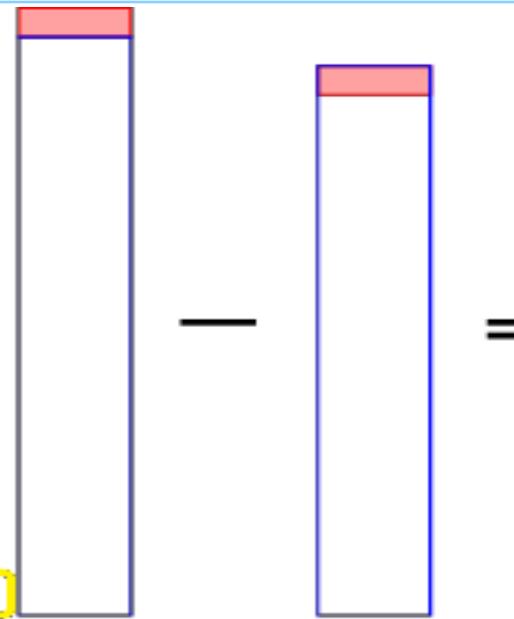


1 Huge relative error for small root $\frac{1}{\gamma}$

* w.r.b. EPS

(14)

We witness cancellation:



▷ Huge error amplification when subtracting two numbers of about the same size.

Fig. 39

C++11-code 1.5.41: Discriminant formula for the real roots of $p(\xi) = \xi^2 + \alpha\xi + \beta$

```

2 //! C++ function computing the zeros of a quadratic polynomial
3 //!  $\xi \rightarrow \xi^2 + \alpha\xi + \beta$  by means of the familiar discriminant
4 //! formula  $\xi_{1,2} = \frac{1}{2}(-\alpha \pm \sqrt{\alpha^2 - 4\beta})$ . However
5 //! this implementation is vulnerable to round-off! The zeros are
6 //! returned in a column vector
7 Vector2d zerosquadpol(double alpha, double beta){
8     Vector2d z;
9     double D = std::pow(alpha, 2) - 4*beta; // discriminant
10    if(D < 0) throw "no real zeros";
11    else{
12        // The famous discriminant formula
13        double wD = std::sqrt(D);
14        z << (-alpha-wD)/2, (-alpha+wD)/2; //
15    }
16    return z;
17 }
```

cancellation here

$$\begin{aligned} \gamma &\gg 1 \\ \alpha &\gg 1, \beta = 1 \\ \Rightarrow D &\approx \alpha \end{aligned}$$

Ex: Difference quotients

$$f'(x) \approx \frac{f(x+h) - f(x)}{h}, h \ll 1$$

cancellation will happen here

Analysis: approximation error $\rightarrow 0$ as $h \rightarrow 0$

C++11-code 1.5.48: Difference quotient approximation of the derivative of exp $\Rightarrow f(x) = e^x$

```

2 //! Difference quotient approximation
3 //! of the derivative of exp
4 void diffq(){
5     double h = 0.1, x = 0.0;
6     for(int i = 1; i <= 16; ++i){
7         double df = (exp(x+h)-exp(x))/h;
8         cout << setprecision(14) << fixed;
9         cout << setw(5) << -i
10            << setw(20) << df-1 << endl;
11         h /= 10;
12     }
13 }
```

$\log_{10}(h)$	relative error
-1	0.05170918075648
-2	0.00501670841679
-3	0.00050016670838
-4	0.00005000166714
-5	0.00000500000696
-6	0.00000049996218
-7	0.00000004943368
-8	-0.0000000607747
-9	0.00000008274037
-10	0.00000008274037
-11	0.00000008274037
-12	0.00008890058234
-13	-0.00079927783736
-14	-0.00079927783736
-15	0.11022302462516
-16	-1.000000000000000

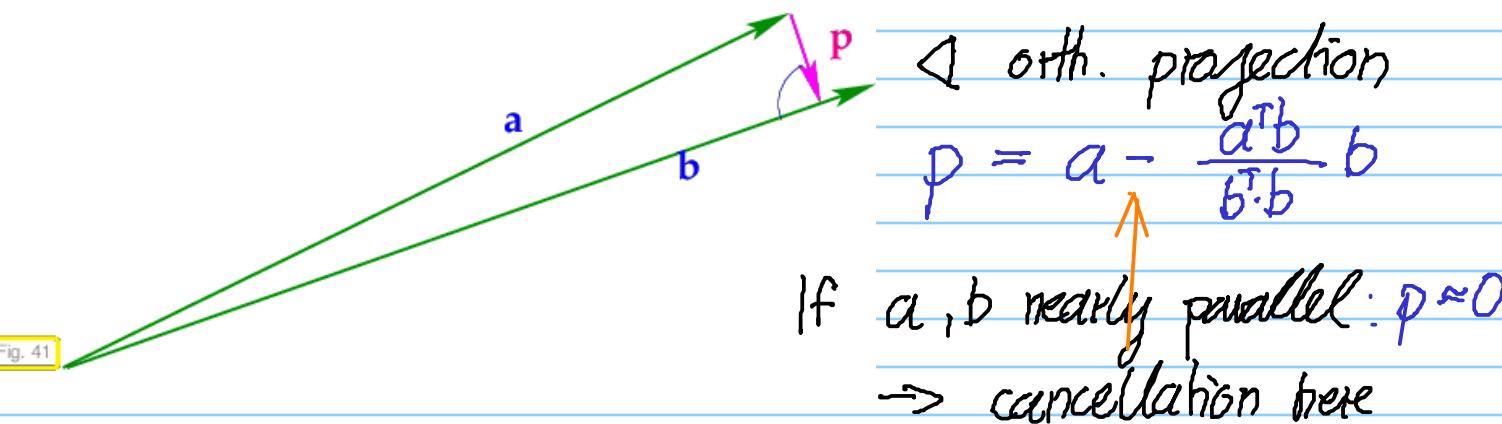
Measured relative errors ▷

We observe an initial decrease of the relative approximation error followed by a steep increase when h drops below 10^{-8} .

$$h_{opt} \approx \sqrt{\text{EPS}}$$

effect of cancellation

Rem: Cancellation & orthogonalization



Avoiding cancellation:

- Computing the roots of a quadratic polynomial:

$$p(\xi) = \xi^2 + \alpha\xi + \beta \quad \text{w/ zeros } \bar{\xi}_1, \bar{\xi}_2$$

$$\Rightarrow \bar{\xi}_1 \cdot \bar{\xi}_2 = \beta$$

Idea: Compute large root $\bar{\xi}_2$ (no problem)

$$\bar{\xi}_1 = \frac{\beta}{\bar{\xi}_2}$$

C++11-code 1.5.54: Stable computation of real root of a quadratic polynomial

```

2  ///! C++ function computing the zeros of a quadratic polynomial
3  ///!  $\xi \rightarrow \xi^2 + \alpha\xi + \beta$  by means of the familiar discriminant
4  ///! formula  $\xi_{1,2} = \frac{1}{2}(-\alpha \pm \sqrt{\alpha^2 - 4\beta})$ .
5  ///! This is a stable implementation based on Vieta's theorem.
6  ///! The zeros are returned in a column vector
7  VectorXd zerosquadpolstab(double alpha, double beta){
8      Vector2d z(2);
9      double D = std::pow(alpha, 2) - 4*beta; // discriminant
10     if(D < 0) throw "no real zeros";
11     else{
12         double wD = std::sqrt(D);
13         // Use discriminant formula only for zero far away from 0
14         // in order to avoid cancellation. For the other zero
15         // use Vieta's formula.
16         if(alpha >= 0){
17             double t = 0.5*(-alpha-wD); //
18             z << t, beta/t;
19         }
20         else{
21             double t = 0.5*(-alpha+wD); //
22             z << beta/t, t;
23         }
24     }
25     return z;
26 }
```

→ no cancellation

- Avoiding cancellation by math. identities

$$\int_0^x \sin t dt = 1 - \cos x = 2 \sin^2(x/2)$$

cancellation for $x \approx 0$

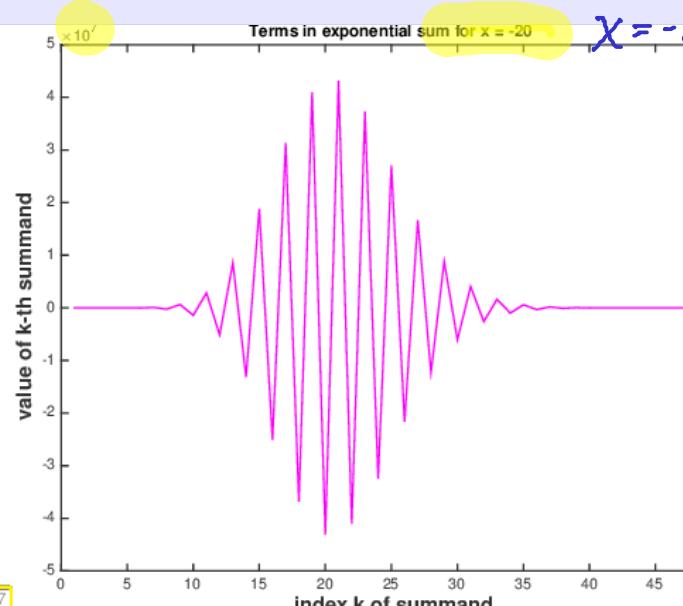
16

Summation of exp-series

C++11-code 1.5.61: Summation of exponential series

```

2 double expeval(double x,
3                 double tol=1e-8){
4     // Initialization
5     double y = 1.0, term = 1.0;
6     long int k = 1;
7     // Termination criterion
8     while(abs(term) > tol*y) {
9         term *= x/k;      // next summand
10        y += term; // Summation
11        ++k;
12    }
13    return y;
14 }
```



$$e^x = \sum_{k=0}^{\infty} \frac{1}{k!} x^k$$

← termination criterion



x	Approximation $\tilde{e}^x(x)$	e^x	$\frac{ e^x - \tilde{e}^x(x) }{e^x}$
-20	6.1475618242e-09	2.0611536224e-09	1.982583033727893
-18	1.5983720359e-08	1.5229979745e-08	0.049490585500089
-16	1.1247503300e-07	1.1253517472e-07	0.000534425951530
-14	8.3154417874e-07	8.3152871910e-07	0.000000018591829627
-12	6.1442105142e-06	6.1442123533e-06	0.0000000299321453
-10	4.5399929604e-05	4.5399929762e-05	0.0000000003501044
-8	3.3546262812e-04	3.3546262790e-04	0.0000000000662004
-6	2.4787521758e-03	2.4787521767e-03	0.0000000000332519
-4	1.8315638879e-02	1.8315638889e-02	0.0000000000530724
-2	1.3533528320e-01	1.3533528324e-01	0.0000000000273603
0	1.0000000000e+00	1.0000000000e+00	0.000000000000000
2	7.3890560954e+00	7.3890560989e+00	0.000000000479969
4	5.4598149928e+01	5.4598150033e+01	0.000000001923058
6	4.0342879295e+02	4.0342879349e+02	0.000000001344248
8	2.9809579808e+03	2.9809579870e+03	0.000000002102584
10	2.2026465748e+04	2.2026465795e+04	0.000000002143799
12	1.6275479114e+05	1.6275479142e+05	0.000000001723845
14	1.2026042798e+06	1.2026042842e+06	0.000000003634135
16	8.8861105010e+06	8.8861105205e+06	0.000000002197990
18	6.5659968911e+07	6.5659969137e+07	0.000000003450972
20	4.8516519307e+08	4.8516519541e+08	0.000000004828737

Remedy : $e^{-x} = \frac{1}{e^x}$

8 digits

(17)

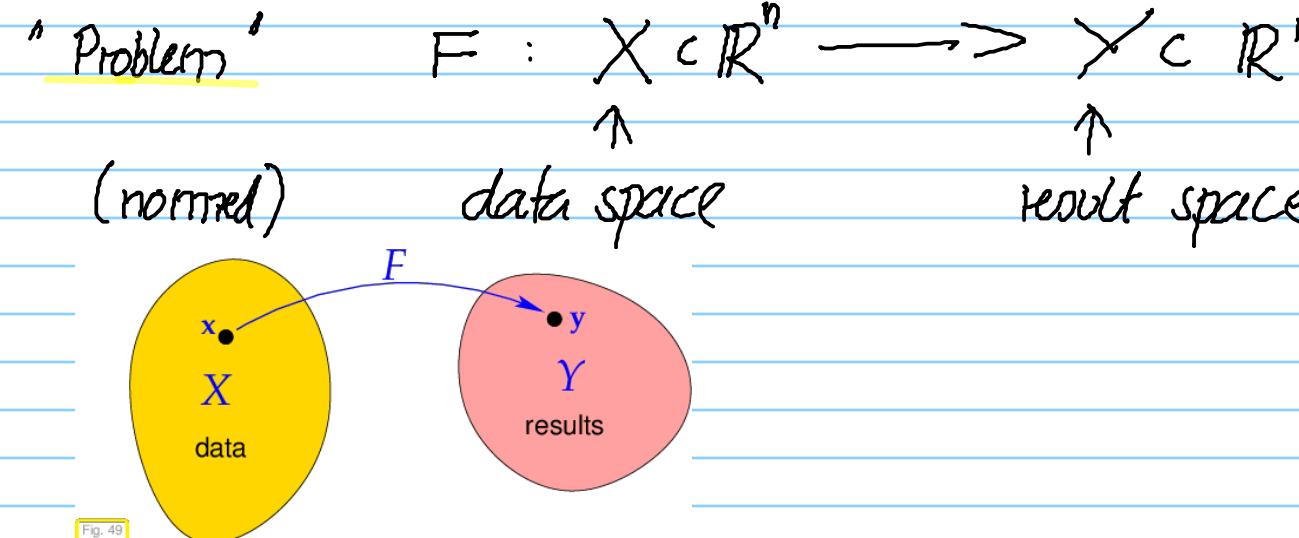
Ex : Defeat cancellation by approximation

$$I(a) := \int_0^a e^{at} dt = \frac{e^{at} - 1}{a} \quad : \text{cancellation for } a \ll 1$$

Idea : $I(a) = \sum_{k=0}^m \underbrace{\frac{1}{(k+1)!} a^k}_{\text{Taylor sum}} + \frac{1}{(m+1)!} a^m \tilde{a}^m$
 $[0 < \tilde{a} < a]$

Choosing m large enough \rightarrow Taylor sum $\approx I(a)$

1.5.5. Numerical Stability



Example: Problem of solving a $n \times n$ (square) linear system of equations (LSE)

$$F : \begin{cases} \mathbb{R}_*^{n \times n} \times \mathbb{R}^n \longrightarrow \mathbb{R}^n \\ (A, \underline{b}) \longrightarrow \underline{x} := A^{-1} \underline{b} \end{cases}$$

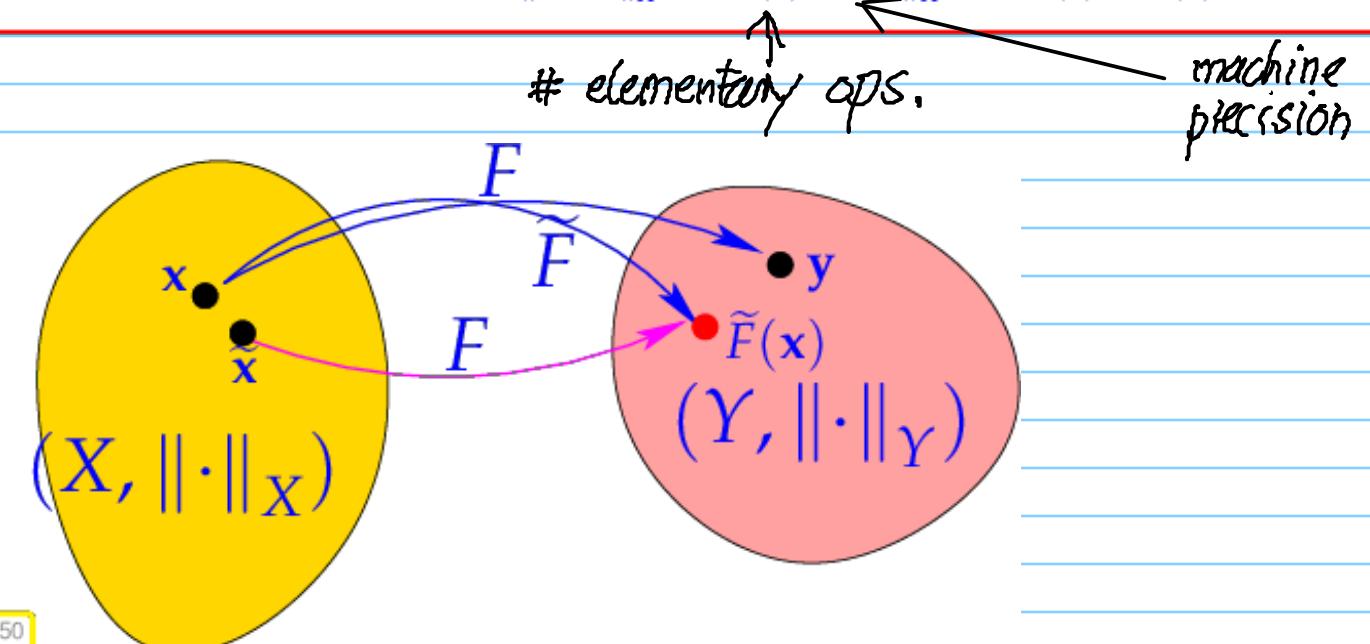
$\mathbb{R}_*^{n \times n} \hat{=} \text{regular } n \times n \text{ matrices}$

Algorithm : $\tilde{F} : \tilde{X} \subset X \longrightarrow \tilde{Y} \subset M^m$
 (work M)

Definition 1.5.82. Stable algorithm = (Good algorithm)

An algorithm \tilde{F} for solving a problem $F : X \mapsto Y$ is numerically stable, if for all $\mathbf{x} \in X$ its result $\tilde{F}(\mathbf{x})$ (possibly affected by roundoff) is the exact result for "slightly perturbed" data:

$$\exists C \approx 1: \forall \mathbf{x} \in X: \exists \tilde{\mathbf{x}} \in X: \|\mathbf{x} - \tilde{\mathbf{x}}\|_X \leq C w(\mathbf{x}) \text{EPS} \|\mathbf{x}\|_X \wedge \tilde{F}(\mathbf{x}) = F(\tilde{\mathbf{x}}).$$



Sloppily speaking, the impact of roundoff (*) on a *stable algorithm* is of the same order of magnitude as the effect of the inevitable perturbations due to rounding the input data.

Stability $\Leftrightarrow \|F(x) - \tilde{F}(x)\| \approx \text{EPS} \|F(x)\|$



Fig. 51

Example : LSE

$\tilde{x} \in \mathbb{R}^n$: perturbed result of solver function

[i.e. $\tilde{x} = \hat{F}(A, b)$]

$$Ax = b \quad [x = F(A, b)]$$

$$\text{to find: } \tilde{A}, \tilde{b} : \quad \tilde{A}\tilde{x} = \tilde{b} \quad [\Leftrightarrow \tilde{x} = F(\tilde{A}, \tilde{b})]$$

$$(i) \quad \tilde{A} = A, \quad \tilde{b} = Ax$$

$$\|\tilde{b} - b\| = \underbrace{\|b - Ax\|}_{!} \leq C \cdot \text{EPS} \cdot \|b\|$$

stable, if this holds for all A, b

\rightarrow prove based on "Axiom"

$$(ii) \quad \tilde{A} = A + \frac{k \tilde{x}^T}{\|\tilde{x}\|^2}, \quad k := b - A\tilde{x}, \quad \tilde{b} = b$$