# Mathematics for chemical engineers

Drahoslava Janovská

## 1. Linear Algebra

**No label** Mandatory material. This will be in writing tests and will be tested during the oral examination.

⭐ Worked examples for practice - optional

⭐ For students who want to know more. This material will not be a part of lectures, it will not be in written tests, and it will not be tested during the oral examination.

## Outline

Matrix equations, inverse of a matrix. Eigenvalues and eigenvectors of a matrix, generalized eigenvectors. Solution of systems of linear algebraic equa

Solution of systems of linear algebraic equations

# Solution of systems of linear algebraic equations

Let us solve a system of linear algebraic equations

$$\mathbf{Ax} = \mathbf{b}\,, \qquad \mathbf{A} \in \mathbb{R}^{m \times n}\,,\ \mathbf{x} \in \mathbb{R}^n\,,\ \mathbf{b} \in \mathbb{R}^m$$

Basically, two types of methods are used:

- direct methods ... after a final number of steps we obtain the solution **x**
- iterative methods ... **x** is gained as a limit of a sequence of iterations $\mathbf{x}_n$ :

$$\mathbf{x} = \lim_{n \longrightarrow \infty} \mathbf{x}_n$$

Here, very important is so called stopping criterion.

**Remak**

We say that the matrix **B** is in upper triangular form (UT–form) iff

$$b_{ii} \neq 0, \quad b_{ij} = 0 \ \ for\ i > j \ (below\ the\ diagonal)\,.$$

Matrix equations, inverse of a matrix. Eigenvalues and eigenvectors of a matrix, generalized eigenvectors. Solution of systems of linear algebraic equa
○●○○○○○○○○○○○○○○○○

Solution of systems of linear algebraic equations

# Direct methods

1. Gauss elimination

$\mathbf{A} \sim \mathbf{B}$ , $\quad$ **B** in upper triangular form (UT–form) . . . forward direction

$$\mathbf{B} = \begin{pmatrix} * & x & x & x \\ 0 & * & x & x \\ 0 & 0 & * & x \end{pmatrix} \uparrow \dots \text{ direction back – the unknowns are evaluated}$$

Theorem of Frobenius

A system of linear algebraic equations $\mathbf{Ax} = \mathbf{b}$ has a solution

$$\Longleftrightarrow$$

$$h(\mathbf{A}) = h(\mathbf{A}|\mathbf{b}) , \text{ i.e.}$$

the rank of the matrix **A** has to be equal to the rank of the extended matrix $(\mathbf{A}|\mathbf{b})$ of the system .

Number of solutions:

- If $h(\mathbf{A}) = h(\mathbf{A}|\mathbf{b}) = n \implies$ the system has just one solution.
- If $h(\mathbf{A}) = h(\mathbf{A}|\mathbf{b}) < n \implies$ the system has an infinite number of solutions and $\dim V_h = n - h(\mathbf{A}) > 0$, where $V_h$ is the space of all solutions of the homogeneous system.

**Matrix equations, inverse of a matrix. Eigenvalues and eigenvectors of a matrix, generalized eigenvectors. Solution of systems of linear algebraic equa**
○○●○○○○○○○○○○○○○○○

Solution of systems of linear algebraic equations

⭐

**Example** Which number $q$ makes this system singular and which right side $t$ gives it infinitely many solutions? Find the solution that has $z = 1$.

$$
\begin{aligned}
x + 4y - 2z &= 1 \\
x + 7y - 6z &= 6 \\
3y + qz &= t.
\end{aligned}
$$

**Solution** Let us transform the extended matrix of the system into UT-form:

$$
\left[\begin{array}{ccc|c}
1 & 4 & -2 & 1 \\
1 & 7 & -6 & 6 \\
0 & 3 & q & t
\end{array}\right] \sim
\left[\begin{array}{ccc|c}
1 & 4 & -2 & 1 \\
0 & 3 & -4 & 5 \\
0 & 3 & q & t
\end{array}\right] \sim
\left[\begin{array}{ccc|c}
1 & 4 & -2 & 1 \\
0 & 3 & -4 & 5 \\
0 & 0 & q+4 & t-5
\end{array}\right]
$$

- $q \neq -4 \ \wedge \ t \neq 5 \ \Rightarrow \ h(\mathbf{A}) = h(\mathbf{A}|\mathbf{b}) = n = 3 \ \Rightarrow$ in this case the system has just one solution:

$$
x = -\frac{17q + 10t + 18}{3(q+4)}, \quad y = \frac{4t + 5q}{3(q+4)}, \quad z = \frac{t-5}{q+4}.
$$

The equation $z = 1$ is fulfilled by all points of the line $t = q + 9$.

Matrix equations, inverse of a matrix. Eigenvalues and eigenvectors of a matrix, generalized eigenvectors. Solution of systems of linear algebraic equa

○○○●○○○○○○○○○○○○

Solution of systems of linear algebraic equations

★

- $q + 4 = 0 \ \wedge \ t - 5 \neq 0 \ \Rightarrow \ q = -4, t \neq 5$. Then $h(\mathbf{A}) = 2$, $h(\mathbf{A}|\mathbf{b}) = 3$ and there is no solution of the system. The number $q = -4$ makes this system singular.

- $q = -4, t = 5$. Then the original system is similar to the system

$$\left[ \begin{array}{ccc|c} 1 & 4 & -2 & 1 \\ 0 & 3 & -4 & 5 \end{array} \right] \ \Rightarrow \ h(\mathbf{A}) = h(\mathbf{A}|\mathbf{b}) = 2, \ n = 3 \ \Rightarrow$$

the system has infinitely many solutions, one unknown has to be a parameter, $z = \alpha, \alpha \in \mathbb{R}$. Then $y = \frac{5}{3} + \frac{4}{3}\alpha$, $x = -\frac{17}{3} - \frac{10}{3}\alpha$.

$$\mathbf{x} = \left[ \begin{array}{c} x \\ y \\ z \end{array} \right] = \frac{1}{3} \left[ \begin{array}{c} -17 \\ 5 \\ 0 \end{array} \right] + \frac{\alpha}{3} \left[ \begin{array}{c} -10 \\ 4 \\ 3 \end{array} \right], \ V_h = \{\mathbf{x} \in \mathbb{R}^3, \mathbf{x} = \left[ \begin{array}{c} x \\ y \\ z \end{array} \right] = \frac{\alpha}{3} \left[ \begin{array}{c} -10 \\ 4 \\ 3 \end{array} \right], \ \alpha \in \mathbb{R} \},$$

$\dim V_h = 1$ and $z = 1 \Leftrightarrow \alpha = 1$. Then $x = -9$, $y = 3$. For $q = -4$ the right hand side $t = 5$ gives infinitely many solutions, $\mathbf{x} = (-9, 3, 1)^\mathsf{T}$ is the solution with $z = 1$

Matrix equations, inverse of a matrix. Eigenvalues and eigenvectors of a matrix, generalized eigenvectors. Solution of systems of linear algebraic equa

○○○○●○○○○○○○○○○○

Solution of systems of linear algebraic equations

# ★ Solution of the system of linear algebraic equations via decompositions

**2.a)** LU–decomposition (Gauss elimination is a particular example of LU–decomposition)

$$\mathbf{A} = \mathbf{LU}$$

**L** ... a lower triangular matrix with ones on the diagonal,
**U** ... an upper triangular matrix, $u_{ii} \neq 0$ (if this is not true, the rows of the matrix has first to be permute by a permutation matrix **P** and then we decompose the matrix **PA**).

$$\mathbf{L} = \begin{pmatrix} 1 & 0 & 0 & 0 \\ x & 1 & 0 & 0 \\ & & \ddots & \ddots \\ x & \dots & x & 1 \end{pmatrix} \qquad \mathbf{U} = \begin{pmatrix} * & x & x & x \\ 0 & * & x & x \\ & & \ddots & \ddots \\ 0 & \dots & 0 & * \end{pmatrix}$$

Because

$$\mathbf{Ax} = \mathbf{b} \iff (\mathbf{LU})\mathbf{x} = \mathbf{b} \iff \mathbf{L}(\mathbf{Ux}) = \mathbf{b},$$

we solve two systems with a triangular matrix each:

At first $\quad \mathbf{Ly} = \mathbf{b}, \quad$ then $\quad \mathbf{Ux} = \mathbf{y}.$

Matrix equations, inverse of a matrix. Eigenvalues and eigenvectors of a matrix, generalized eigenvectors. Solution of systems of linear algebraic equa
○○○○○●○○○○○○○○○○
Solution of systems of linear algebraic equations

★

**2.b)** QR–decomposition

$$\mathbf{A} = \mathbf{QR}$$

$\mathbf{Q}$ ... an orthogonal matrix: $\quad \mathbf{QQ}^\mathrm{T} = \mathbf{E} \quad \Longleftrightarrow \quad \mathbf{Q}^{-1} = \mathbf{Q}^\mathrm{T}$
$\mathbf{R}$ ... an upper triangular matrix.

We multiply the equation $(\mathbf{QR})\mathbf{x} = \mathbf{b}$ from the left by the matrix $\mathbf{Q}^\mathrm{T}$.
We obtain

$$\underbrace{\mathbf{Q}^\mathrm{T}\mathbf{Q}}_{\mathbf{E}}\mathbf{R}\mathbf{x} = \mathbf{Q}^\mathrm{T}\mathbf{b} \qquad \mathbf{Rx} = \mathbf{Q}^\mathrm{T}\mathbf{b}$$

and we have again a system with a triangular matrix.

As we have already noted, the LU–decomposition need not exist.
QR–decomposition exists always.

LU–decomposition is advantageous namely if we have to solve many
linear systems with the same matrix. Then we will compute the matrices
$\mathbf{L}$ and $\mathbf{U}$ only for the first system.

Matrix equations, inverse of a matrix. Eigenvalues and eigenvectors of a matrix, generalized eigenvectors. Solution of systems of linear algebraic equa
ooooooo●oooooooooo

Solution of systems of linear algebraic equations

# Other direct methods

**3.** Cramer's rule – it is ideal for regular matrices $2 \times 2$, for larger systems it is unusable .

**4.** Application of the inverse $\mathbf{A}^{-1}$ of the matrix $\mathbf{A}$ ($\mathbf{E}$ is the identity matrix)

$$\underbrace{\mathbf{A}\mathbf{x} = \mathbf{b}}, \quad \mathbf{A} \text{ regular} \quad \Longrightarrow \quad \exists \mathbf{A}^{-1} : \mathbf{A}\mathbf{A}^{-1} = \mathbf{A}^{-1}\mathbf{A} = \mathbf{E}.$$

the simplest matrix equation

We multiply the equation $\mathbf{A}\mathbf{x} = \mathbf{b}$ from the left by the matrix $\mathbf{A}^{-1}$:

$$\underbrace{\mathbf{A}^{-1}\mathbf{A}}_{\mathbf{E}}\mathbf{x} = \mathbf{A}^{-1}\mathbf{b} \quad \Longrightarrow \quad \mathbf{x} = \mathbf{A}^{-1}\mathbf{b}$$

This method is very unstable it can't be recommended for practical numerical computations. It is useful in proofs of theoretical results.

Matrix equations, inverse of a matrix. Eigenvalues and eigenvectors of a matrix, generalized eigenvectors. Solution of systems of linear algebraic equa...
○○○○○○○●○○○○○○○○

Solution of systems of linear algebraic equations

# Iterative methods

**5.** Jacobi and Gauss–Seidel methods
The idea:

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

$\mathbf{D} = \operatorname{diag}(a_{ii})$ a diagonal matrix

$-\mathbf{L}$ a sharp lower triangular matrix

$-\mathbf{U}$ a sharp upper triangular matrix

$$\mathbf{A} = \begin{pmatrix} a_{11} & x & x & x \\ x & a_{22} & x & x \\ x & x & a_{33} & x \\ x & x & x & a_{44} \end{pmatrix}$$

Let $a_{ii} \neq 0$, $i = 1, \dots, n$, then the matrix $\mathbf{L} + \mathbf{U}$ has all diagonal entries equal to zero.

$$\mathbf{Ax} = \mathbf{b} \quad \Longleftrightarrow \quad (\mathbf{D} - \mathbf{L} - \mathbf{U})\mathbf{x} = \mathbf{b}$$

$$\mathbf{D} \underbrace{\mathbf{x}}_{(k+1)-\text{st iteration}} = \mathbf{b} + (\mathbf{L} + \mathbf{U}) \underbrace{\mathbf{x}}_{k-\text{th iteration}} \quad \Longrightarrow \quad \underbrace{\mathbf{x}^{(k+1)} = \mathbf{D}^{-1}(\mathbf{b} + (\mathbf{L} + \mathbf{U})\mathbf{x}^{(k)})}_{\text{the Jacobi method}}$$

Matrix equations, inverse of a matrix. Eigenvalues and eigenvectors of a matrix, generalized eigenvectors. Solution of systems of linear algebraic equa
○○○○○○○○●○○○○○○○○

Solution of systems of linear algebraic equations

## The Gauss–Seidel method

Let the matrix **A** be again written as

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

$$(\mathbf{D} - \mathbf{L} - \mathbf{U})\mathbf{x} = \mathbf{b} \quad \Longleftrightarrow \quad (\mathbf{D} - \mathbf{L})\underbrace{\mathbf{x}}_{(k+1)\text{st iteration}} = \mathbf{b} + \mathbf{U}\underbrace{\mathbf{x}}_{k-\text{th iteration}} \quad \Longrightarrow$$

the Gauss–Seidel method :

$$\mathbf{x}^{(k+1)} = (\mathbf{D} - \mathbf{L})^{-1}(\mathbf{b} + \mathbf{U}\mathbf{x}^{(k)}) \, .$$

Matrix equations, inverse of a matrix. Eigenvalues and eigenvectors of a matrix, generalized eigenvectors. Solution of systems of linear algebraic equa
○○○○○○○○○○●○○○○○○○

Matrix inversion

## Matrix inversion

Let **A** be a square matrix $n \times n$. We say that a matrix **B** is an inverse matrix to the matrix **A**, iff $\mathbf{AB} = \mathbf{BA} = \mathbf{E}$. The matrix inverse to **A** is usually denoted as $\mathbf{A}^{-1}$, i.e.,

$$\mathbf{AA}^{-1} = \mathbf{A}^{-1}\mathbf{A} = \mathbf{E}.$$

The inverse matrix $\mathbf{A}^{-1}$ to the matrix **A** exists $\iff$ **A** is nonsingular $(\det\mathbf{A} \neq \mathbf{0})$.

- Let $\mathbf{A} = (a_{ij})_{n \times n}$, $\det \mathbf{A} \neq 0$ (**A** is a regular matrix). Then

$$\mathbf{A}^{-1} = \frac{1}{\det \mathbf{A}} \begin{pmatrix} \mathbf{A}_{11}, & \ldots & \mathbf{A}_{1n} \\ \mathbf{A}_{21}, & \ldots & \mathbf{A}_{2n} \\ \vdots & & \vdots \\ \mathbf{A}_{n1}, & \ldots & \mathbf{A}_{nn} \end{pmatrix}^{\mathrm{T}}, \quad \text{where } \mathbf{A}_{ij} = (-1)^{i+j}\, \mathbf{M}_{ij},$$

$\mathbf{A}_{ij}$ is a cofactor of $a_{ij}$, $\mathbf{M}_{ij}$ is a minor that belongs to $a_{ij}$.

- Gauss–Jordan method
  In this case we don't need to know at the beginning of the computation that the matrix **A** is regular. We use equivalent operations to transform

$$(\mathbf{A}|\mathbf{E}) \sim (\mathbf{E}|\mathbf{A}^{-1}).$$

Matrix equations, inverse of a matrix. Eigenvalues and eigenvectors of a matrix, generalized eigenvectors. Solution of systems of linear algebraic equa

○○○○○○○○○○○●○○○○○

Matrix equations

## Matrix equations

The simplest matrix equation

$$\mathbf{A}\mathbf{x} = \mathbf{b}, \quad \mathbf{A} \text{ square, regular}, \ n \times n \quad \Longrightarrow \quad \exists \mathbf{A}^{-1}$$

We multiply the equation by $\mathbf{A}^{-1}$ from the left and obtain

$$\underbrace{\mathbf{A}^{-1}\mathbf{A}}_{\mathbf{E}}\mathbf{x} = \mathbf{A}^{-1}\mathbf{b} \quad \Longrightarrow \quad \mathbf{x} = \mathbf{A}^{-1}\mathbf{b}$$

If I would multiply the equation by $\mathbf{A}^{-1}$ from the right:

$$\underbrace{\mathbf{A}}_{n \times n} \underbrace{\mathbf{x}}_{n \times 1} \underbrace{\mathbf{A}^{-1}}_{n \times n} = \underbrace{\mathbf{b}}_{n \times 1} \underbrace{\mathbf{A}^{-1}}_{n \times n} \quad \dots \quad \text{can't multiply}$$

**Example:**

$$\mathbf{X}\mathbf{A} - \mathbf{E} = 2\mathbf{X} + \mathbf{A}, \qquad \mathbf{A} = \begin{pmatrix} 1 & 3 \\ 0 & 2 \end{pmatrix}$$

$$\mathbf{X}\mathbf{A} - 2\mathbf{X} = \mathbf{A} + \mathbf{E} \quad \Longrightarrow \quad \mathbf{X}(\mathbf{A} - 2\mathbf{E}) = \mathbf{A} + \mathbf{E}, \quad \mathbf{A} - 2\mathbf{E} = \begin{pmatrix} -1 & 3 \\ 0 & 0 \end{pmatrix}$$

$\det(\mathbf{A} - 2\mathbf{E}) = 0 \quad \Longrightarrow \mathbf{A} - 2\mathbf{E}$ doesn't have an inversion.

The matrix equation has no solution.

Matrix equations, inverse of a matrix. Eigenvalues and eigenvectors of a matrix, generalized eigenvectors. Solution of systems of linear algebraic equa

○○○○○○○○○○○○●○○○○○

**Matrix equations**

⭐

**Example**

$$\mathbf{A}\,\mathbf{X} + \mathbf{X} = \mathbf{A} - \mathbf{E}, \qquad \mathbf{A} = \begin{pmatrix} 1 & 1 & 1 \\ 1 & 1 & 0 \\ 1 & 0 & 0 \end{pmatrix}$$

$$(\mathbf{A} + \mathbf{E})\,\mathbf{X} = \mathbf{A} - \mathbf{E} \quad \Longrightarrow \quad \mathbf{X} = (\mathbf{A} + \mathbf{E})^{-1}\,(\mathbf{A} - \mathbf{E})$$

$$\mathbf{A} + \mathbf{E} = \begin{pmatrix} 2 & 1 & 1 \\ 1 & 2 & 0 \\ 1 & 0 & 1 \end{pmatrix}, \qquad \det(\mathbf{A} + \mathbf{E}) = 1 \neq 0 \quad \Longrightarrow \quad \exists\,(\mathbf{A} + \mathbf{E})^{-1}$$

$$(\mathbf{A} + \mathbf{E}|\mathbf{E}) = \left( \begin{array}{ccc|ccc} 2 & 1 & 1 & 1 & 0 & 0 \\ 1 & 2 & 0 & 0 & 1 & 0 \\ 1 & 0 & 1 & 0 & 0 & 1 \end{array} \right) \sim \left( \begin{array}{ccc|ccc} 2 & 1 & 1 & 1 & 0 & 0 \\ 0 & -3 & 1 & 1 & -2 & 0 \\ 0 & 1 & -1 & 1 & 0 & -2 \end{array} \right) \sim$$

$$\sim \left( \begin{array}{ccc|ccc} 6 & 0 & 4 & 4 & -2 & 0 \\ 0 & -3 & 1 & 1 & -2 & 0 \\ 0 & 0 & -2 & 4 & -2 & -6 \end{array} \right) \sim \left( \begin{array}{ccc|ccc} 6 & 0 & 0 & 12 & -6 & -12 \\ 0 & -6 & 0 & 6 & -6 & -6 \\ 0 & 0 & -2 & 4 & -2 & -6 \end{array} \right) \sim$$

$$\sim \left( \begin{array}{ccc|ccc} 1 & 0 & 0 & 2 & -1 & -2 \\ 0 & 1 & 0 & -1 & 1 & 1 \\ 0 & 0 & 1 & -2 & 1 & 3 \end{array} \right) = (\mathbf{E}|(\mathbf{A} + \mathbf{E})^{-1})$$

$$\mathbf{X} = (\mathbf{A} + \mathbf{E})^{-1}\,(\mathbf{A} - \mathbf{E}) = \begin{pmatrix} -3 & 2 & 4 \\ 2 & -1 & -2 \\ 4 & -2 & -5 \end{pmatrix}.$$

Matrix equations, inverse of a matrix. Eigenvalues and eigenvectors of a matrix, generalized eigenvectors. Solution of systems of linear algebraic equa

○○○○○○○○○○○○○●○○○

Eigenvalues and eigenvectors of the matrix

# Eigenvalues and eigenvectors of the matrix

**Definition** A number $\lambda$ (real or complex) is called the eigenvalue of the real or complex matrix **A** if it satisfies for a nonzero vector **x** the equation

$$\mathbf{A}\mathbf{x} = \lambda\,\mathbf{x}\,,$$

$\mathbf{x} \neq \mathbf{0}$ ... eigenvector of the matrix **A** that corresponds to the eigenvalue $\lambda$.

The set of all eigenvalues of the matrix **A** ... spectrum of the matrix **A**

$\mathbf{A}\mathbf{x} = \lambda\mathbf{x}$ ... the matrix equation for the unknown eigenvector **x**

$\mathbf{A}\mathbf{x} - \lambda\mathbf{x} = \mathbf{0}\,,\quad (\mathbf{A} - \lambda\mathbf{E})\mathbf{x} = \mathbf{0}\,,\ \mathbf{x} \neq \mathbf{0}\,,\quad \mathbf{A} - \lambda\mathbf{E}$ must be singular $\implies$

$\underbrace{\det(\mathbf{A} - \lambda\mathbf{E}) = 0}$ ... characteristic equation of the matrix **A**

characteristic polynomial of the matrix **A** = polynomial of degree $n$:

$$P(\lambda) = \det(\mathbf{A} - \lambda\mathbf{E}) = (-1)^n(\lambda^n + p_1\lambda^{n-1} + p_2\lambda^{n-2} + \cdots + p_n)\,,\quad \text{w}here$$

$$-p_1 = a_{11} + a_{22} + \cdots + a_{nn} = \text{ the trace of the matrix } \mathbf{A}$$
$$p_n = (-1)^n\det\mathbf{A}$$

Be careful! Eigenvalues of a real matrix may be imaginary.

Matrix equations, inverse of a matrix. Eigenvalues and eigenvectors of a matrix, generalized eigenvectors. Solution of systems of linear algebraic equa

○○○○○○○○○○○○○○○●○○

Eigenvalues and eigenvectors of the matrix

⭐

**Example**  Let us compute eigenvalues and corresponding eigenvectors of
the matrix $\mathbf{A} = \begin{pmatrix} 2 & 1 \\ -5 & 0 \end{pmatrix}$ .

$$\mathbf{A} - \lambda\mathbf{E} = \begin{pmatrix} 2-\lambda & 1 \\ -5 & -\lambda \end{pmatrix} \quad \Longrightarrow \quad \det(\mathbf{A} - \lambda\mathbf{E}) = \lambda^2 - 2\lambda + 5$$

Characteristic equation  $\lambda^2 - 2\lambda + 5 = 0 \implies \lambda_1 = 1 + 2i, \quad \lambda_2 = 1 - 2i$ ,
eigenvalues of the matrix $\mathbf{A}$ are the complex conjugate numbers.
Let us calculate eigenvector $\mathbf{x}_1$ that corresponds to the eigenvalue
$\lambda_1 = 1 + 2i$ , i.e., we have to solve the system with the singular matrix :

$$(\mathbf{A} - \lambda_1\mathbf{E})\mathbf{x}_1 = \mathbf{0}, \ \mathbf{x}_1 = (h_1, h_2)^{\mathrm{T}} \neq \mathbf{0}, \quad \text{i.e., the system}$$
$$\begin{pmatrix} 1-2i & 1 \\ -5 & -1-2i \end{pmatrix} \begin{pmatrix} h_1 \\ h_2 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

This system has infinity many solutions, but we need only one eigenvector.
Let us choose for example  $h_1 = 1$ , then  $h_2 = -1 + 2i$ . Because the
eigenvalues are complex conjugate, the eigenvectors are also complex
conjugate. We obtain

$$\mathbf{x}_1 = \begin{pmatrix} 1 \\ -1+2i \end{pmatrix}, \qquad \mathbf{x}_2 = \begin{pmatrix} 1 \\ -1-2i \end{pmatrix}.$$

Matrix equations, inverse of a matrix. Eigenvalues and eigenvectors of a matrix, generalized eigenvectors. Solution of systems of linear algebraic equa

○○○○○○○○○○○○○○○●○

Eigenvalues and eigenvectors of the matrix

## An estimate of the spectral radius

An estimate of the spectral radius – the Gershgorin theorem
Let $\mathbf{A} = (a_{jk})$ be a square $n \times n$ matrix. Let us denote

$$K_j = \{\mu \in \mathbb{C}, \ |\mu - a_{jj}| \leq \sum_{\substack{k=1 \\ k \neq j}}^{n} |a_{jk}|\} \quad \text{is a circle with the center } S_j \text{ and radius } r_j,$$

$$S_j = a_{jj}, \ r_j = \sum_{\substack{k=1 \\ k \neq j}}^{n} |a_{jk}|, \text{ i.e., the radius } r_j \text{ of the circle } K_j \text{ is equal to the sum}$$

of absolute values of non diagonal entries in $j$−th row. Then all eigenvalues

of the matrix $\mathbf{A}$ are located in the union of all circles, i.e., in $\bigcup_{j=1}^{n} K_j$.

Numerical methods for computation of eigenvalues are based on LU or QR
decomposition of the matrix $\mathbf{A}$.

Matrix equations, inverse of a matrix. Eigenvalues and eigenvectors of a matrix, generalized eigenvectors. Solution of systems of linear algebraic equa

○○○○○○○○○○○○○○○●

**Eigenvalues and eigenvectors of the matrix**

**Example**

$$\mathbf{A} = \begin{pmatrix} 1 & -1 & 1 & -1 \\ 0 & 1 & -1 & 0 \\ -1 & 0 & 0 & -1 \\ 0 & 1 & 2 & -2 \end{pmatrix}, \quad \begin{array}{rcl} S_1 &=& 1 \\ S_2 &=& 1 \\ S_3 &=& 0 \\ S_4 &=& -2 \end{array} \quad \begin{array}{rcl} r_1 &=& 3 \\ r_2 &=& 1 \\ r_3 &=& 2 \\ r_4 &=& 3 \end{array}$$



Eigenvalues are

$$1.126575852 \quad \pm \quad 0.7768133722\,i, \quad -1.126575852 \quad \pm \quad 1.391009448\,i\,.$$

All eigenvalues are located in the set $M = K_1 \cup K_2 \cup K_3 \cup K_4$ .

Matrix equations, inverse of a matrix. Eigenvalues and eigenvectors of a matrix, generalized eigenvectors. Solution of systems of linear algebraic equa

○○○○○○○○○○○○○○○○○○○

**The Givens matrices of plane rotations**

# ★ The Givens matrices of plane rotations

The matrix $\mathbf{G}_{pq} \in \mathbb{R}^{n \times n}$, $p < q$, of the form

$$
\mathbf{G}_{p,q} =
\begin{pmatrix}
1 & & & & & & & & & & \\
& \ddots & & & & & & & & & \\
& & 1 & & & & & & & & \\
& & & c & \cdots\cdots\cdots & s & & & & & \\
& & & \vdots & 1 & & \vdots & & & & \\
& & & \vdots & & \ddots & \vdots & & & & \\
& & & \vdots & & & 1 & \vdots & & & \\
& & & -s & \cdots\cdots\cdots & & c & & & \\
& & & & & & & 1 & & & \\
& & & & & & & & \ddots & \\
& & & & & & & & & 1
\end{pmatrix}
\begin{matrix}
\\ \\ \\ \longleftarrow p \\ \\ \\ \\ \longleftarrow q \\ \\ \\
\end{matrix} \quad ,
$$

$$
\begin{matrix}
\uparrow & & \uparrow \\
p & & q
\end{matrix}
$$

where $s = \sin \phi$, $c = \cos \phi$, $\phi$ real, is called the Givens matrix of the plane rotation.

The matrix $\mathbf{G}_{pq}$ performs the rotation of $\mathbb{R}^n$ around the point $0 \in \mathbb{R}^n$ about the angle $\phi$ in $p$-$q$ plane.

Matrix equations, inverse of a matrix. Eigenvalues and eigenvectors of a matrix, generalized eigenvectors. Solution of systems of linear algebraic equa

○○○○○○○○○○○○○○○○○○○

**The Givens matrices of plane rotations**

# ★ Elementary rotation matrix

The idea:

$$\mathbf{G}_{12} = \begin{pmatrix} \cos\varphi & \sin\varphi \\ -\sin\varphi & \cos\varphi \end{pmatrix}, \; \mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \implies \mathbf{G}_{12}\mathbf{x} = \begin{pmatrix} x_1\cos\varphi + x_2\sin\varphi \\ -x_1\sin\varphi + x_2\cos\varphi \end{pmatrix}.$$

By a suitable choice of $\varphi$ we can gain that one of the component of the vector **x** vanishes. For example, let us choose $\varphi$ in such a way that the second component will vanish:

$$-x_1\sin\varphi + x_2\cos\varphi = 0. \quad \text{Now, we apply} \quad \cos\varphi = \sqrt{1 - \sin^2\varphi}\,.$$

Then

$$x_1^2\sin^2\varphi = x_2^2(1 - \sin^2\varphi) \implies \sin^2\varphi = \frac{x_2^2}{x_1^2 + x_2^2} \quad \text{pro } (x_1, x_2) \neq (0,0)\,.$$

We have

$$\sin\varphi = \frac{|x_2|}{\sqrt{x_1^2 + x_2^2}}; \; \text{similarly} \; \cos\varphi = \frac{|x_1|}{\sqrt{x_1^2 + x_2^2}}\,.$$

(it is not necessary to know the angel $\varphi$, we need only $\sin\varphi$ and $\cos\varphi$).

Matrix equations, inverse of a matrix. Eigenvalues and eigenvectors of a matrix, generalized eigenvectors. Solution of systems of linear algebraic equa
○○○○○○○○○○○○○○○○○○○○○

The Givens matrices of plane rotations

# ★ Zeroing one component of the vector

$\mathbf{G}_{1,2} \ldots$ elementary matrix of the rotation.

$\mathbf{G}_{1,2}$ is the orthogonal matrix:

$$\mathbf{G}_{1,2}^{\mathrm{T}} \mathbf{G}_{1,2} = \mathbf{G}_{1,2} \mathbf{G}_{1,2}^{\mathrm{T}} = \begin{pmatrix} \cos\varphi & \sin\varphi \\ -\sin\varphi & \cos\varphi \end{pmatrix} \begin{pmatrix} \cos\varphi & -\sin\varphi \\ \sin\varphi & \cos\varphi \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} = \mathbf{E}$$

Let us set

$$\sin\varphi = \frac{|x_2|}{\sqrt{x_1^2 + x_2^2}} \quad \text{a} \quad \cos\varphi = \frac{|x_1|}{\sqrt{x_1^2 + x_2^2}} .$$

Then

$$\mathbf{G}_{12}\mathbf{x} = \begin{pmatrix} x_1 \cos\varphi + x_2 \sin\varphi \\ -x_1 \sin\varphi + x_2 \cos\varphi \end{pmatrix} = \begin{pmatrix} \dfrac{x_1^2}{\sqrt{x_1^2 + x_2^2}} + \dfrac{x_2^2}{\sqrt{x_1^2 + x_2^2}} \\ \dfrac{-x_1 x_2}{\sqrt{x_1^2 + x_2^2}} + \dfrac{x_1 x_2}{\sqrt{x_1^2 + x_2^2}} \end{pmatrix} = \begin{pmatrix} \|\mathbf{x}\| \\ 0 \end{pmatrix} .$$

Matrix equations, inverse of a matrix. Eigenvalues and eigenvectors of a matrix, generalized eigenvectors. Solution of systems of linear algebraic equa
○○○○○○○○○○○○○○○○○

**The Givens matrices of plane rotations**

★

**Example**   Let $\mathbf{x} = (3, 4)^{\mathrm{T}}$. Let us set $\sin \phi = \frac{4}{5}$, $\cos \phi = \frac{3}{5}$. Then

$$\mathbf{G}_{12}\mathbf{x} = \frac{1}{5} \begin{pmatrix} 3 & 4 \\ -4 & 3 \end{pmatrix} \begin{pmatrix} 3 \\ 4 \end{pmatrix} = (5, 0)^{\mathrm{T}}.$$

If we multiply the vector $\mathbf{x}^{\mathrm{T}}$ by the matrix $\mathbf{G}_{pq}^{\mathrm{T}}$ from the right, we have

$$\mathbf{x}^{\mathrm{T}}\mathbf{G}_{12}^{\mathrm{T}} = \frac{1}{5}( 3, \ 4 ) \begin{pmatrix} 3 & -4 \\ 4 & 3 \end{pmatrix} = ( 5, \ 0 ).$$

If we put $\sin \phi = \frac{3}{5}$, $\cos \phi = -\frac{4}{5}$, we obtain
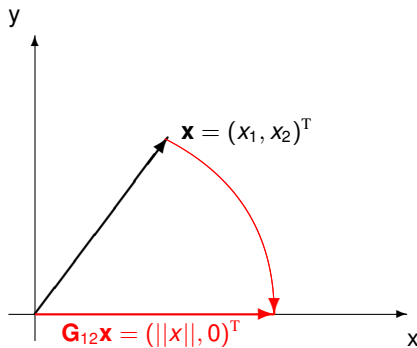
$$\mathbf{G}_{12}\mathbf{x} = \frac{1}{5} \begin{pmatrix} -4 & 3 \\ -3 & -4 \end{pmatrix} \begin{pmatrix} 3 \\ 4 \end{pmatrix} = (0, -5)^{\mathrm{T}},$$

and similarly,

$$\mathbf{x}^{\mathrm{T}}\mathbf{G}_{12}^{\mathrm{T}} = \frac{1}{5}( 3, \ 4 ) \begin{pmatrix} -4 & -3 \\ 3 & -4 \end{pmatrix} = ( 0, \ -5 ).$$

Matrix equations, inverse of a matrix. Eigenvalues and eigenvectors of a matrix, generalized eigenvectors. Solution of systems of linear algebraic equations
○○○○○○○○○○○○○○○○○○○

**The Givens matrices of plane rotations**

$$\mathbf{G}_{12}\mathbf{x} = \begin{pmatrix} \|\mathbf{x}\| \\ 0 \end{pmatrix}$$

Matrix equations, inverse of a matrix. Eigenvalues and eigenvectors of a matrix, generalized eigenvectors. Solution of systems of linear algebraic equations

ooooooooooooooooo

**The Givens matrices of plane rotations**

# ★ Rotation of a vector in $\mathbb{R}^3$

The matrix of rotation **G** for a vector $\mathbf{x} \in \mathbb{R}^3$?

The idea: we substitute $\mathbf{G}_{1,2}$ into the identity matrix $\mathbf{E}_{3\times 3}$ - three possibilities:

- For a suitable choice of $\varphi$

$$\widetilde{\mathbf{G}}_{12} = \begin{pmatrix} \cos\varphi & \sin\varphi & 0 \\ -\sin\varphi & \cos\varphi & 0 \\ 0 & 0 & 1 \end{pmatrix} \implies \widetilde{\mathbf{G}}_{12} \cdot \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} \widetilde{x_1} \\ 0 \\ x_3 \end{pmatrix}$$

- For a suitable choice of $\varphi$

$$\widetilde{\mathbf{G}}_{13} = \begin{pmatrix} \cos\varphi & 0 & \sin\varphi \\ 0 & 1 & 0 \\ -\sin\varphi & 0 & \cos\varphi \end{pmatrix} \implies \widetilde{\mathbf{G}}_{13} \cdot \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} \widetilde{x_1} \\ x_2 \\ 0 \end{pmatrix}$$

- For a suitable choice of $\varphi$

$$\widetilde{\mathbf{G}}_{23} \begin{pmatrix} 1 & 0 & 0 \\ 0 & \cos\varphi & \sin\varphi \\ 0 & -\sin\varphi & \cos\varphi \end{pmatrix} \implies \widetilde{\mathbf{G}}_{23} \cdot \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} x_1 \\ \widetilde{x_2} \\ 0 \end{pmatrix}$$

If we apply the Givens plane rotation to any vector, only one of the components of this vector will vanish.

Matrix equations, inverse of a matrix. Eigenvalues and eigenvectors of a matrix, generalized eigenvectors. Solution of systems of linear algebraic equa
○○○○○○○○○○○○○○○○○

The Givens matrices of plane rotations

$$\mathbf{G} := \widetilde{\widetilde{\mathbf{G}_{13}}}\widetilde{\mathbf{G}_{12}} \implies \mathbf{G} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \widetilde{\widetilde{\mathbf{G}_{13}}} \begin{pmatrix} \widetilde{x}_1 \\ 0 \\ x_3 \end{pmatrix} = \begin{pmatrix} \widetilde{\widetilde{x}}_1 \\ 0 \\ 0 \end{pmatrix}$$

We applied two rotation matrices and as a result we obtained the vector with two zero components of the vector **x**.

$\mathbf{G}_{12}$ is an orthogonal matrix $\implies \mathbf{G}_{12}^{-1} = \mathbf{G}_{12}^{T}$

$$\mathbf{x}^{T}\mathbf{G}_{12}^{-1} = (x_1, x_2) \cdot \mathbf{G}_{12}^{-1} = (x_1, x_2) \cdot \begin{pmatrix} \cos\varphi & -\sin\varphi \\ sin\varphi & \cos\varphi \end{pmatrix} =$$

$$= \left( \frac{(x_1)^2}{\sqrt{x_1^2 + x_2^2}} + \frac{(x_2)^2}{\sqrt{x_1^2 + x_2^2}}, -\frac{x_1 x_2}{\sqrt{x_1^2 + x_2^2}} + \frac{x_2 x_1}{\sqrt{x_1^2 + x_2^2}} \right) = (||x||, 0)$$

If we multiply $\mathbf{x}^{T}$ from the right by the matrix $\mathbf{G}_{12}^{-1}$ the second component of the vector $\mathbf{x}^{T}$ will vanish.

Matrix equations, inverse of a matrix. Eigenvalues and eigenvectors of a matrix, generalized eigenvectors. Solution of systems of linear algebraic equa
○○○○○○○○○○○○○○○○

**The Givens matrices of plane rotations**

★

The Givens matrix of the plane rotation $\mathbf{G}_{pq} \in \mathbb{R}^{n \times n}$ is always the orthogonal matrix

$$\mathbf{G}_{pq}^{\mathrm{T}} \mathbf{G}_{pq} = \mathbf{G}_{pq} \mathbf{G}_{pq}^{\mathrm{T}} = \mathbf{E} \, .$$

It is different from the identity matrix $\mathbf{E}$ only in entries in positions

$$(p, p), \ (p, q), \ (q, p) \ \text{a} \ (q, q) \, .$$

If we multiply any matrix $\mathbf{A}$ by the matrix $\mathbf{G}_{pq}$ from the left, only $p$th and $q$th column of the matrix $\mathbf{A}$ will change, if we multiply any matrix $\mathbf{A}$ by the matrix $\mathbf{G}_{pq}^{\mathrm{T}}$ from the right, only $p$th and $q$th row of the matrix $\mathbf{A}$ will change.

Matrix equations, inverse of a matrix. Eigenvalues and eigenvectors of a matrix, generalized eigenvectors. Solution of systems of linear algebraic equa

○○○○○○○○○○○○○○○○○○

**The Givens matrices of plane rotations**

## ★ The order of zeroing the entries

By a suitable choice of the Givens matrices one can transform a given matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ into the similar upper triangular matrix. The order of zeroing the entries is very important, because we don't want the entries that have already vanished to become nonzero again.

Let us construct a sequence of the matrices $\mathbf{G}_{pq}$ $(p < q)$ in such a way that $q$th component of the given vector will vanish and we multiply the matrix $\mathbf{A}$ from the left by these matrices in the following order

$$\begin{matrix} \mathbf{G}_{12} & \mathbf{G}_{13} & \ldots & \mathbf{G}_{1n} \\ & \mathbf{G}_{23} & \ldots & \mathbf{G}_{2n} \\ & & \ddots & \\ & & & \mathbf{G}_{n-1,n} \end{matrix} \quad .$$

After $k$ multiplications, $k \leq \frac{1}{2} n(n-1)$ (we don't apply the Givens rotation to that entry of the matrix which is already zero) we obtain the similar matrix $\mathbf{R}$ in the upper triangular form. The matrix

$$\mathbf{G} = \mathbf{G}_{n-1,n} \mathbf{G}_{n-2,n} \ldots \mathbf{G}_{13} \mathbf{G}_{12}$$

is orthogonal,

$$\mathbf{G}\,\mathbf{A} = \mathbf{R} \quad \Longrightarrow \quad \mathbf{A} = \mathbf{G}^{\mathrm{T}}\,\mathbf{R}.$$

Matrix equations, inverse of a matrix. Eigenvalues and eigenvectors of a matrix, generalized eigenvectors. Solution of systems of linear algebraic equa
○○○○○○○○○○○○○○○○○○

**The Givens matrices of plane rotations**

★

Let us draw the schema of the Givens method applied to the matrix $\mathbf{A} \in \mathbb{R}^{4 \times 4}$ in order to transform it into upper triangular form. ($+ \ldots$ the entries that have not been changed by the transformation, $* \ldots$ elements that changed):

$$\mathbf{A} = \begin{pmatrix} + & + & + & + \\ + & + & + & + \\ + & + & + & + \\ + & + & + & + \end{pmatrix} \xrightarrow{\mathbf{G}_{12}} \begin{pmatrix} * & * & * & * \\ 0 & * & * & * \\ + & + & + & + \\ + & + & + & + \end{pmatrix} \xrightarrow{\mathbf{G}_{13}} \begin{pmatrix} * & * & * & * \\ 0 & + & + & + \\ 0 & * & * & * \\ + & + & + & + \end{pmatrix} \longrightarrow$$

$$\xrightarrow{\mathbf{G}_{14}} \begin{pmatrix} * & * & * & * \\ 0 & + & + & + \\ 0 & + & + & + \\ 0 & * & * & * \end{pmatrix} \xrightarrow{\mathbf{G}_{23}} \begin{pmatrix} + & + & + & + \\ 0 & * & * & * \\ 0 & 0 & * & * \\ 0 & + & + & + \end{pmatrix} \xrightarrow{\mathbf{G}_{24}} \begin{pmatrix} + & + & + & + \\ 0 & * & * & * \\ 0 & 0 & + & + \\ 0 & 0 & * & * \end{pmatrix} \xrightarrow{\mathbf{G}_{34}} \begin{pmatrix} + & + & + & + \\ 0 & + & + & + \\ 0 & 0 & * & * \\ 0 & 0 & 0 & * \end{pmatrix} = \mathbf{R}.$$

Disadvantage: The Givens matrix sets always to zero only one element in the given matrix.

Advantage (namely if you work with sparse matrices): If we apply the Givens rotation to a given matrix then only two columns or rows will be changed. Others entries will not change.

Matrix equations, inverse of a matrix. Eigenvalues and eigenvectors of a matrix, generalized eigenvectors. Solution of systems of linear algebraic equa

○○○○○○○○○○○○○○○○

**Householder's matrices of the reflection**

## ★ Householder's matrices of the reflection

Householder's matrix of the reflection ... the matrix that is able to set to zero more entries of the given vector in once.

$$\mathbf{H_v} = \mathbf{E} - 2\,\frac{\mathbf{v}\,\mathbf{v}^{\mathrm{T}}}{||\mathbf{v}||^2}\,, \quad \mathbf{v} \in \mathbb{R}^n,\ \mathbf{v} \neq \mathbf{0}\,.$$

$\mathbf{x} \in \mathbb{R}^n \implies \mathbf{H}\mathbf{x} := \mathbf{H_v}\mathbf{x}$ is a vector symmetric with the given vector $\mathbf{x}$ by the manifold $\varrho$, that is orthogonal to the vector $\mathbf{v}$.

$$\mathbf{H}^{\mathrm{T}} = (\mathbf{E} - 2\,\frac{\mathbf{v}\,\mathbf{v}^{\mathrm{T}}}{||\mathbf{v}||^2})^{\mathrm{T}} = \mathbf{H} \implies \mathbf{H}\ \text{is symmetric}$$

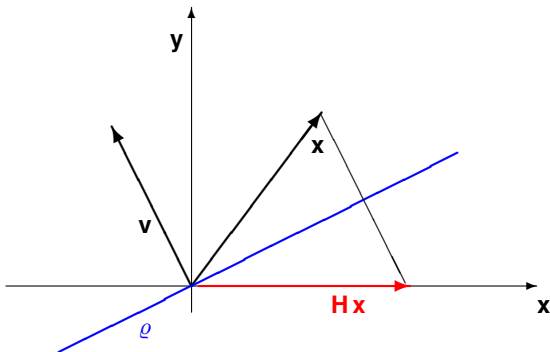$$\mathbf{H}^{\mathrm{T}}\mathbf{H} = \mathbf{H}^2 = \mathbf{E} \implies \mathbf{H}\ \text{is orthogonal}\,.$$

**Remark:** $\mathbf{H}\mathbf{x} = \mathbf{x} - \mathbf{v} \implies$ for $\mathbf{v} = \mathbf{x} - ||\mathbf{x}||\mathbf{e}_1$ is $\mathbf{H}\mathbf{x} = ||\mathbf{x}||\mathbf{e}_1$
the vector that has all components except the first one equal to zero.

**Householder's matrices of the reflection**

**Example**

$$\mathbf{x} = (3, 4)^{\mathrm{T}}, \quad \mathbf{v} = \mathbf{x} - \|\mathbf{x}\|\mathbf{e}_1 = (-2, 4)^{\mathrm{T}},$$

$$\mathbf{H} = \frac{1}{5} \begin{pmatrix} 3 & 4 \\ 4 & -3 \end{pmatrix}, \quad \mathbf{H}\mathbf{x} = (5, 0)^{\mathrm{T}}.$$

Matrix equations, inverse of a matrix. Eigenvalues and eigenvectors of a matrix, generalized eigenvectors. Solution of systems of linear algebraic equa

oooooooooooooooooo

**Householder's matrices of the reflection**

## ★ Solution of the system of linear algebraic equations by the Householder method

$$\mathbf{A}\mathbf{x} = \mathbf{b}\,, \quad \mathbf{A} \in \mathbb{R}^{n \times n}\,.$$

We find $n-1$ Householder's matrices $\mathbf{H}_1,\ \mathbf{H}_2,\ \ldots,\mathbf{H}_{n-1}$, such that

$$\mathbf{H}\,\mathbf{A} := \mathbf{H}_{n-1} \cdot \cdots \cdot \mathbf{H}_2 \cdot \mathbf{H}_1 \cdot \mathbf{A} = \mathbf{R}\,,$$

$\mathbf{R}$ is an upper triangular matrix. Then we solve the system with the triangular matrix $\mathbf{R}$:

$$\mathbf{H}\,\mathbf{A} = \mathbf{R}\,, \quad \mathbf{H} \quad \text{is orthogonal, i.e.} \quad \mathbf{H}^{\mathrm{T}} = \mathbf{H}^{-1} \quad \Longrightarrow \quad \mathbf{A} = \mathbf{H}^{\mathrm{T}}\,\mathbf{R}\,.$$

$$\mathbf{H}\,\mathbf{A}\mathbf{x} = \mathbf{H}\mathbf{b} \quad \Longrightarrow \quad \mathbf{R}\mathbf{x} = \mathbf{H}\mathbf{b}$$

Matrix equations, inverse of a matrix. Eigenvalues and eigenvectors of a matrix, generalized eigenvectors. Solution of systems of linear algebraic equa
○○○○○○○○○○○○○○○○○

Singular values

# Singular values of a rectangular matrix

$\mathbf{A} \in \mathbb{R}^{m \times n}$, $m \neq n$ ... any rectangular matrix. We can't define eigenvalues for rectangular matrices, but ...

$$\mathbf{A}^{\mathrm{T}}\mathbf{A} \in \mathbb{R}^{n \times n} \implies \mathbf{A}^{\mathrm{T}}\mathbf{A} \text{ is a square matrix}$$

$$\left(\mathbf{A}^{\mathrm{T}}\mathbf{A}\right)^{\mathrm{T}} = \mathbf{A}^{\mathrm{T}}\mathbf{A} \implies \mathbf{A}^{\mathrm{T}}\mathbf{A} \text{ is symmetric}$$

$$\mathbf{x}^{\mathrm{T}}\left(\mathbf{A}^{\mathrm{T}}\mathbf{A}\right)\mathbf{x} = (\mathbf{A}\mathbf{x})^{\mathrm{T}}\mathbf{A}\mathbf{x} = \|\mathbf{A}\mathbf{x}\|^2 \geq 0 \,\forall\, \mathbf{x} \in \mathbb{R}^n \implies \mathbf{A}^{\mathrm{T}}\mathbf{A} \text{ is positive semidefinite}$$

Eigenvalues $\lambda_1, \lambda_2, \ldots, \lambda_n$ of the matrix $\mathbf{A}^{\mathrm{T}}\mathbf{A}$ are real, nonnegative. We may write them as $\lambda_k = \sigma_k^2$, $\sigma_k \geq 0$, $k = 1, \ldots, n$. The numbers

$\sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_n \geq 0$ are called singular values of the matrix $\mathbf{A}$.

For the largest and the smallest singular value of the rectangular matrix $\mathbf{A}$, it holds:

$$\sigma_1 = \max_{\mathbf{0} \neq \mathbf{x} \in \mathbb{R}^n} \frac{\|\mathbf{A}\mathbf{x}\|}{\|\mathbf{x}\|}, \quad \sigma_n = \min_{\mathbf{0} \neq \mathbf{x} \in \mathbb{R}^n} \frac{\|\mathbf{A}\mathbf{x}\|}{\|\mathbf{x}\|}.$$

Matrix equations, inverse of a matrix. Eigenvalues and eigenvectors of a matrix, generalized eigenvectors. Solution of systems of linear algebraic equ

Singular value decomposition

## Singular value decomposition

**Theorem** Let $\mathbf{A} \in \mathbb{R}^{m \times n}$ be any matrix. Then

- there exist two orthogonal matrices $\mathbf{U} \in \mathbb{R}^{m \times m}$ and $\mathbf{V} \in \mathbb{R}^{n \times n}$ such that the $m \times n$ matrix $\mathbf{S} = \mathbf{U}^{\mathrm{T}}\mathbf{A}\mathbf{V}$ has a "diagonal"form

$$\mathbf{S} = \left( \begin{array}{cc} \mathbf{D} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{array} \right)$$

$$\mathbf{D} = \text{diag}\,(\sigma_1, \sigma_2, \ldots, \sigma_r), \quad \sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_r > 0,$$

where $\sigma_1, \sigma_2, \ldots, \sigma_r$ are the nonzero singular values of the matrix $\mathbf{A}$ and $r$ is the rank of the matrix $\mathbf{A}$;

- the nonzero singular values of the matrix $\mathbf{A}^{\mathrm{T}}$ are also the numbers $\sigma_1, \sigma_2, \ldots, \sigma_r$.

The decomposition $\mathbf{A} = \mathbf{U}\mathbf{S}\mathbf{V}^{\mathrm{T}}$ ... singular value decomposition of the matrix $\mathbf{A}$.

**Remark** $\mathbf{S} = \mathbf{U}^{\mathrm{T}}\mathbf{A}\mathbf{V}$

- the columns of the matrix $\mathbf{U}$ ... $m$ orthonormal eigenvectors of the symmetric $m \times m$ matrix $\mathbf{A}\mathbf{A}^{\mathrm{T}}$,
- the columns of the matrix $\mathbf{V}$ ... $n$ orthonormal eigenvectors of the symmetric $n \times n$ matrix $\mathbf{A}^{\mathrm{T}}\mathbf{A}$.

Matrix equations, inverse of a matrix. Eigenvalues and eigenvectors of a matrix, generalized eigenvectors. Solution of systems of linear algebraic equa
○○○○○○○○○○○○○○○○

Singular value decomposition

★

**Example**    Compute the singular value decomposition of the matrix

$$\mathbf{A} = \left[ \begin{array}{ccc} 1 & 1 & 0 \\ 0 & 1 & 1 \end{array} \right].$$

**Solution**    We want to compute the decomposition $\mathbf{A} = \mathbf{U}\mathbf{S}\mathbf{V}^{\mathrm{T}}$, where $\mathbf{U}$ is $2 \times 2$ orthogonal matrix, $\mathbf{S}$ has singular values on the diagonal, $\mathbf{V}$ is $3 \times 3$ orthogonal matrix.

$$\mathbf{A}^{\mathrm{T}}\mathbf{A} = \left[ \begin{array}{ccc} 1 & 1 & 0 \\ 1 & 2 & 1 \\ 0 & 1 & 1 \end{array} \right]$$

Eigenvalues of $\mathbf{A}^{\mathrm{T}}\mathbf{A}$ :    $\lambda_1 = 3, \lambda_2 = 1, \lambda_3 = 0$

Singular values of $\mathbf{A}$ :    $\sigma_1 = \sqrt{3}, \sigma_2 = 1, \sigma_3 = 0$

Corresponding eigenvectors:

$$\lambda_1 = 3 : \quad (\mathbf{A}^{\mathrm{T}}\mathbf{A} - 3\mathbf{E}) = \left[ \begin{array}{ccc} -2 & 1 & 0 \\ 1 & -1 & 1 \\ 0 & 1 & -2 \end{array} \right] \sim \left[ \begin{array}{ccc} -2 & 1 & 0 \\ 0 & -1 & 2 \end{array} \right]$$

$$\widetilde{\mathbf{v}}_1 = (1, 2, 1)^{\mathrm{T}}, \ \|\widetilde{\mathbf{v}}_1\| = \sqrt{6} \quad \Rightarrow \quad \mathbf{v}_1 = \frac{1}{\sqrt{6}}(1, 2, 1)^{\mathrm{T}}.$$

Matrix equations, inverse of a matrix. Eigenvalues and eigenvectors of a matrix, generalized eigenvectors. Solution of systems of linear algebraic equa
○○○○○○○○○○○○○○○○○○○○

Singular value decomposition

$$\lambda_2 = 1: \quad (\mathbf{A}^{\mathrm{T}}\mathbf{A} - \mathbf{E}) = \begin{bmatrix} 0 & 1 & 0 \\ 1 & 1 & 1 \\ 0 & 1 & 0 \end{bmatrix} \sim \begin{bmatrix} 1 & 1 & 1 \\ 0 & 1 & 0 \end{bmatrix}$$

$$\widetilde{\mathbf{v}}_2 = (1, 0, -1)^{\mathrm{T}}, \ \|\widetilde{\mathbf{v}_2}\| = \sqrt{2} \quad \Rightarrow \ \mathbf{v}_2 = \frac{1}{\sqrt{2}}(1, 0, -1)^{\mathrm{T}}.$$

$$\lambda_3 = 0: \quad (\mathbf{A}^{\mathrm{T}}\mathbf{A}) = \begin{bmatrix} 1 & 1 & 0 \\ 1 & 2 & 1 \\ 0 & 1 & 1 \end{bmatrix} \sim \begin{bmatrix} 1 & 1 & 0 \\ 0 & 1 & 1 \end{bmatrix}$$

$$\widetilde{\mathbf{v}}_3 = (1, -1, 1)^{\mathrm{T}}, \ \|\widetilde{\mathbf{v}_3}\| = \sqrt{3} \quad \Rightarrow \ \mathbf{v}_3 = \frac{1}{\sqrt{3}}(1, -1, 1)^{\mathrm{T}}.$$

We obtained

$$\mathbf{V} = \begin{bmatrix} \dfrac{1}{\sqrt{6}} & \dfrac{1}{\sqrt{2}} & \dfrac{1}{\sqrt{3}} \\ \dfrac{2}{\sqrt{6}} & 0 & -\dfrac{1}{\sqrt{3}} \\ \dfrac{1}{\sqrt{6}} & -\dfrac{1}{\sqrt{2}} & \dfrac{1}{\sqrt{3}} \end{bmatrix} \quad \text{and} \quad \mathbf{S} = \begin{bmatrix} \sqrt{3} & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix}$$

Matrix equations, inverse of a matrix. Eigenvalues and eigenvectors of a matrix, generalized eigenvectors. Solution of systems of linear algebraic equa
○○○○○○○○○○○○○○○○○○

**Singular value decomposition**

⭐

We multiply the equation $\mathbf{A} = \mathbf{U}\mathbf{S}\mathbf{V}^{\mathrm{T}}$ by the matrix $\mathbf{V}$ from the right and obtain

$$\mathbf{A}\mathbf{V} = \mathbf{U}\mathbf{S} \quad \Rightarrow \quad \mathbf{A}\mathbf{v}_i = \sigma_i\mathbf{u}_i\,,$$

and we can compute the columns of the matrix $\mathbf{U}$ by formula

$$\mathbf{u}_i = \frac{1}{\sigma_i}\mathbf{A}\mathbf{v}_i\,, \quad \text{i.e.,} \quad \mathbf{u}_1 = \frac{1}{\sigma_1}\mathbf{A}\mathbf{v}_1 = \frac{1}{\sqrt{2}}\begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \mathbf{u}_2 = \frac{1}{\sigma_2}\mathbf{A}\mathbf{v}_2 = \frac{1}{\sqrt{2}}\begin{bmatrix} 1 \\ -1 \end{bmatrix}$$

Thus,

$$\mathbf{U} = \frac{1}{\sqrt{2}}\begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}\,.$$

We can check that the equation $\mathbf{A} = \mathbf{U}\mathbf{S}\mathbf{V}^{\mathrm{T}}$ really holds.

Matrix equations, inverse of a matrix. Eigenvalues and eigenvectors of a matrix, generalized eigenvectors. Solution of systems of linear algebraic equa

○○○○○○○○○○○○○○○○○○○

**Singular value decomposition**

# Practical computation of the singular value decomposition

Let $\mathbf{A} \in \mathbb{R}^{m \times n}$. Is there any connection between spectral decompositions of the matrices $\mathbf{A}^{\mathrm{T}}\mathbf{A}$ and $\mathbf{A}\mathbf{A}^{\mathrm{T}}$?

Let us perform the spectral decomposition of the matrix $\mathbf{A}^{\mathrm{T}}\mathbf{A}$, i.e., we compute eigenvalues $\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_r > 0$ of the matrix $\mathbf{A}^{\mathrm{T}}\mathbf{A}$ and corresponding orthonormal eigenvectors $v_1, v_2, \ldots, v_n$.

Then the orthonormal eigenvectors $u_1, u_2, \ldots u_m$ corresponding to the same nonzero eigenvalues of the matrix $\mathbf{A}\mathbf{A}^{\mathrm{T}}$ are obtained by a simple formula

$$u_j = \frac{Av_j}{\sigma_j}, \quad j = 1, \ldots, m.$$

It is not necessary to perform spectral decomposition of the matrix $\mathbf{A}\mathbf{A}^{\mathrm{T}}$.

Matrix equations, inverse of a matrix. Eigenvalues and eigenvectors of a matrix, generalized eigenvectors. Solution of systems of linear algebraic equa
○○○○○○○○○○○○○○○○○

**Singular value decomposition**

## ★ Numerical computation

$$\mathbf{A} \in \mathbb{R}^{m \times n}, \quad m \geq n \text{ (w.l.o.g.)}$$

Golub–Reinsch algorithm:   two steps

- bidiagonalization of the matrix **A** by applying the Householder matrices of the reflection:

$$\mathbf{A} \longrightarrow \mathbf{J}^{(0)} = \begin{pmatrix} \mathbf{J}_0 \\ \mathbf{0} \end{pmatrix}, \quad \mathbf{J}_0 = \begin{pmatrix} x & x & & & 0 \\ & x & x & & \\ & & \ddots & \ddots & \\ & & & \ddots & x \\ 0 & & & & x \end{pmatrix}.$$

After $n$ reduction steps we obtain an upper bidiagonal $m \times n$ matrix $\mathbf{J}^{(0)}$,

$$\mathbf{J}^{(0)} = \mathbf{P}_n \mathbf{P}_{n-1} \ldots \mathbf{P}_1 \mathbf{A} \mathbf{Q}_1 \mathbf{Q}_2 \ldots \mathbf{Q}_{n-2},$$

$\mathbf{P}_k$, $\mathbf{Q}_k$ are the Householder matrices of the reflection.

Matrix equations, inverse of a matrix. Eigenvalues and eigenvectors of a matrix, generalized eigenvectors. Solution of systems of linear algebraic equa
○○○○○○○○○○○○○○○○
**Singular value decomposition**

★

$$\mathbf{Q} := \mathbf{Q}_1\mathbf{Q}_2\cdots\mathbf{Q}_{n-2}, \quad \mathbf{P} := \mathbf{P}_1\mathbf{P}_2\cdots\mathbf{P}_n, \quad \mathbf{P} \text{ and } \mathbf{Q} \ldots \text{ orthogonal},$$

$$\mathbf{J}^{(0)} = \mathbf{P}^{\mathrm{T}}\mathbf{A}\mathbf{Q}, \quad (\mathbf{J}^{(0)})^{\mathrm{T}}\mathbf{J}^{(0)} = \mathbf{J}_0^{\mathrm{T}}\mathbf{J}_0 = \mathbf{Q}^{\mathrm{T}}\mathbf{A}^{\mathrm{T}}\mathbf{A}\mathbf{Q}.$$

The matrices $\mathbf{J}_0$ and $\mathbf{A}$ are similar, i.e. they have the same singular values. Now, we have to perform the singular value decomposition of the bidiagonal matrix $\mathbf{J}_0$.

- We diagonalize $\mathbf{J}_0$ by using a particular variant of the QR method with shifts based on applying of the sequence of the Givens rotation matrices

$$\mathbf{J}_0 \longrightarrow \mathbf{J}_1 \longrightarrow \ldots \longrightarrow \mathbf{D}, \quad \text{where} \quad \mathbf{D} \text{ is diagonal}, \quad \mathbf{J}_{k+1} = \mathbf{S}_k^{\mathrm{T}}\mathbf{J}_k\mathbf{T}_k,$$

$\mathbf{S}_k$ and $\mathbf{T}_k$ are orthogonal matrices. We will choose the matrices $\mathbf{T}_k$ in such a way that the sequence of the tridiagonal matrices $\mathbf{M}_k = \mathbf{J}_k^{\mathrm{T}}\mathbf{J}_k$ converge to a diagonal matrix. The matrices $\mathbf{S}_k$ are chosen such that all matrices $\mathbf{J}_k$ are in a bidiagonal form.

The method is very quick and numerically stable. The details are out of the aim of the course. You can read more about the method in Wilkinson J. H., Reinsch C. or in Golub G., Van Loan Ch.

Matrix equations, inverse of a matrix. Eigenvalues and eigenvectors of a matrix, generalized eigenvectors. Solution of systems of linear algebraic equa
○○○○○○○○○○○○○○○○
**Overdetermined system of linear equations**

## Overdetermined system of linear equations

$\mathbf{A} \in \mathbb{R}^{m \times n}, \quad m \geq n, \quad \mathbf{b} \in \mathbb{R}^m, \quad ? \, \mathbf{x} \in \mathbb{R}^n : \quad \mathbf{Ax} = \mathbf{b}$.

Our system has more equations then unknowns, we say it is overdetermined.
The system has a solution only if the right hand side $\mathbf{b}$ is an element of the
column vector $\mathcal{R}(\mathbf{A})$ of the matrix $\mathbf{A}$:

$$\mathbf{b} \in \mathcal{R}(\mathbf{A}) = \{\mathbf{y} \in \mathbb{R}^m, \ \exists \mathbf{x} \in \mathbb{R}^n : \mathbf{Ax} = \mathbf{y}\} \subset \mathbb{R}^m,$$

i.e. $\mathbf{b}$ must be a linear combination of the columns of the matrix $\mathbf{A}$. The
components of the computed vector $\mathbf{x} = (x_1, x_2, \ldots, x_n)^{\mathrm{T}}$ are the coefficients
of this linear combination:

$$x_1 \begin{pmatrix} a_{11} \\ a_{21} \\ \vdots \\ a_{m1} \end{pmatrix} + x_2 \begin{pmatrix} a_{12} \\ a_{22} \\ \vdots \\ a_{m2} \end{pmatrix} + \cdots + x_n \begin{pmatrix} a_{1n} \\ a_{2n} \\ \vdots \\ a_{mn} \end{pmatrix} = \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_m \end{pmatrix}.$$
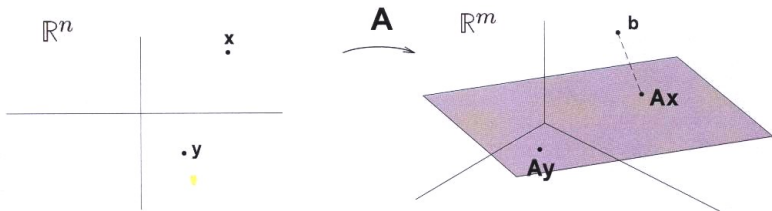
Usually, $\mathbf{b} \notin \mathcal{R}(\mathbf{A})$ and insted of seeking the exact solution, we have to look
for such a solution that is the "closest one" to the exact solution.

Matrix equations, inverse of a matrix. Eigenvalues and eigenvectors of a matrix, generalized eigenvectors. Solution of systems of linear algebraic equa
○○○○○○○○○○○○○○○○○○○

**Overdetermined system of linear equations**

# Problem formulation

Let $\quad \mathbf{A} \in \mathbb{R}^{m \times n}, \quad m \geq n, \quad \mathbf{b} \in \mathbb{R}^m$.

We are looking for a solution of the system $\mathbf{Ax} = \mathbf{b}$ in sense of the least squares, i.e., we are looking for

$$\mathbf{x} \in \mathbb{R}^n, \quad \mathbf{x} \in \arg \min_{\mathbf{y} \in \mathbb{R}^n} ||\mathbf{Ay} - \mathbf{b}||.$$



The least square solution of the system

Matrix equations, inverse of a matrix. Eigenvalues and eigenvectors of a matrix, generalized eigenvectors. Solution of systems of linear algebraic equa
○○○○○○○○○○○○○○○○○○

**Overdetermined system of linear equations**

# Normal equations

$\|\mathbf{x}\| = (\mathbf{x}^{\mathrm{T}}\mathbf{x})^{\frac{1}{2}}$ ... the norm of the vector $\mathbf{x}$

$E = \|\mathbf{Ax} - \mathbf{b}\|$ ... error of the computation

We are looking for such a point $\widetilde{\mathbf{x}} \in \mathbb{R}^n$, for which the error $E$ is minimal.

Geometrically:  the error $E$ is the distance of the points $\mathbf{Ax}$ and $\mathbf{b}$. This distance is the smallest one if the vector $\mathbf{A}\widetilde{\mathbf{x}}$ is an orthogonal projection of the vector $\mathbf{b}$ on the space $\mathcal{R}(\mathbf{A})$

$\implies$  the error vector $\mathbf{A}\widetilde{\mathbf{x}} - \mathbf{b}$ has to be orthogonal to the space $\mathcal{R}(\mathbf{A})$,

i.e., for any $\mathbf{x} \in \mathbb{R}^n$ the vector $\mathbf{Ax} \in \mathcal{R}(\mathbf{A})$ has to be orthogonal to the vector $\mathbf{A}\widetilde{\mathbf{x}} - \mathbf{b}$:

$$(\mathbf{Ax})^{\mathrm{T}}(\mathbf{A}\widetilde{\mathbf{x}} - \mathbf{b}) = 0 \quad \Leftrightarrow \quad \mathbf{x}^{\mathrm{T}}(\mathbf{A}^{\mathrm{T}}\mathbf{A}\widetilde{\mathbf{x}} - \mathbf{A}^{\mathrm{T}}\mathbf{b}) = 0 \qquad \forall\, \mathbf{x} \in \mathbb{R}^n.$$

The last equation can be fulfilled if and only if the vector $\widetilde{\mathbf{x}}$ solves the so called system of normal equations

$$\mathbf{A}^{\mathrm{T}}\mathbf{A}\widetilde{\mathbf{x}} = \mathbf{A}^{\mathrm{T}}\mathbf{b}\,.$$

Matrix equations, inverse of a matrix. Eigenvalues and eigenvectors of a matrix, generalized eigenvectors. Solution of systems of linear algebraic equ
○○○○○○○○○○○○○○○○○

**Overdetermined system of linear equations**

**Theorem**    The vector $\mathbf{x} \in \mathbb{R}^n$ is a least square solution of the system $\mathbf{Ax} = \mathbf{b}$ if and only if $\mathbf{x}$ solves the system of normal equations.
Moreover, the least square problem has just one solution if and only if the rank $\mathrm{h}(\mathbf{A})$ of the matrix $\mathbf{A}$ is maximal, i.e. $\mathrm{h}(\mathbf{A}) = n$. In this case we say that the matrix $\mathbf{A}$ has the full rank.

**Remark**    If $\mathbf{A} \in \mathbb{R}^{m \times n}$, $m \geq n$, then

$$\mathrm{h}(\mathbf{A}) = n \qquad \Longleftrightarrow \qquad \det(\mathbf{A}^{\mathrm{T}}\mathbf{A}) \neq 0.$$

The system $\mathbf{Ax} = \mathbf{b}$ has a unique least square solution if and only if the matrix $\mathbf{A}^{\mathrm{T}}\mathbf{A}$ is regular.

Then from the normal equations we obtain

$$\widetilde{\mathbf{x}} = (\mathbf{A}^{\mathrm{T}}\mathbf{A})^{-1}\mathbf{A}^{\mathrm{T}}\mathbf{b}$$

and if the matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ has orthonormal columns, i.e., $\mathbf{A}^{\mathrm{T}}\mathbf{A} = \mathbf{E}_n$, where $\mathbf{E}_n$ is $n \times n$ identity matrix, then

$$\widetilde{\mathbf{x}} = \mathbf{A}^{\mathrm{T}}\mathbf{b}\,.$$

Matrix equations, inverse of a matrix. Eigenvalues and eigenvectors of a matrix, generalized eigenvectors. Solution of systems of linear algebraic equa
ooooooooooooooooo

Solution of the normal equations

## ★ Theory: The solution of the normal equations

$$\mathbf{A}^{\mathrm{T}}\mathbf{A}\mathbf{x} = \mathbf{A}^{\mathrm{T}}\mathbf{b}, \quad \mathbf{c} := \mathbf{A}^{\mathrm{T}}\mathbf{b}, \quad \mathbf{A}^{\mathrm{T}}\mathbf{A} \ \text{ nonsingular } \implies \mathbf{x} = (\mathbf{A}^{\mathrm{T}}\mathbf{A})^{-1}\mathbf{c}$$

The spectral analysis of the matrix $\mathbf{A}^{\mathrm{T}}\mathbf{A} \in \mathbb{R}^{n \times n}$ :

eigenvalues $\lambda_i > 0$, eigenvectors $\mathbf{v}_i \in \mathbb{R}^n$, $\{\mathbf{v}_1, \ldots, \mathbf{v}_n\} \ldots$ base of $\mathbb{R}^n \implies$

$$\mathbf{x} = \sum_{i=1}^{n} \alpha_i \mathbf{v}_i, \qquad \mathbf{c} = \sum_{i=1}^{n} \gamma_i \mathbf{v}_i$$

$$\left. \begin{array}{rcl} \mathbf{A}^{\mathrm{T}}\mathbf{A}(\sum_{i=1}^{n} \alpha_i \mathbf{v}_i) & = & \sum_{i=1}^{n} \gamma_i \mathbf{v}_i \\[2mm] \sum_{i=1}^{n} \alpha_i \lambda_i \mathbf{v}_i & = & \sum_{i=1}^{n} \gamma_i \mathbf{v}_i \end{array} \right\} \implies \alpha_i = \frac{1}{\lambda_i}\gamma_i, \quad i = 1, \ldots, n.$$

How interprete the last equation?

Matrix equations, inverse of a matrix. Eigenvalues and eigenvectors of a matrix, generalized eigenvectors. Solution of systems of linear algebraic equa
○○○○○○○○○○○○○○○○
Solution of the normal equations

Let eigenvalues $\lambda_i > 0$, $i = 1, \ldots, n$, of the matrix $\mathbf{A}^{\mathrm{T}}\mathbf{A}$ are such that $\lambda_n$ is much smaller then other eigenvalues:
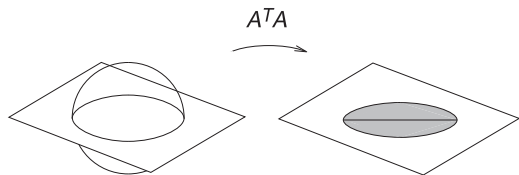
$$\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_{n-1} >> \lambda_n > 0.$$

The matrix $\mathbf{A}^{\mathrm{T}}\mathbf{A}$ maps the sphere in $\mathbb{R}^n$ with the diameter equal to 1 into an ellipsoid with axes in the directions of the eigenvectors $\mathbf{v}_i$. The length of the axis in the direction $\mathbf{v}_n$ is much smaller then lengths of the other axes. It means that the mapping $\mathbf{A}^{\mathrm{T}}\mathbf{A}$ maps any vector of the length 1 into a vector with a negligible $n$th component and the ellipsoid lies in fact in $\mathbb{R}^{n-1}$.

If we solve the normal equations we have to apply the inversion mapping $(\mathbf{A}^{\mathrm{T}}\mathbf{A})^{-1}$. This mapping has the same eigenvectors, but eigenvalues are $\frac{1}{\lambda_i}$:

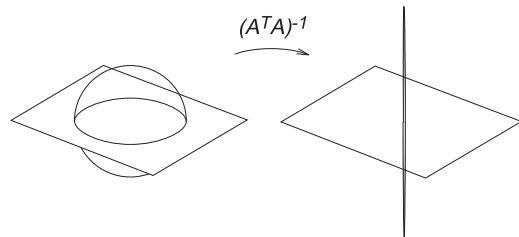$$(\mathbf{A}^{\mathrm{T}}\mathbf{A})^{-1}\mathbf{v}_i = \frac{1}{\lambda_i}\mathbf{v}_i.$$

Because $\frac{1}{\lambda_n}$ is much larger then others $\frac{1}{\lambda_i}$, the corresponding ellipsoid will be in fact onedimensional.

Matrix equations, inverse of a matrix. Eigenvalues and eigenvectors of a matrix, generalized eigenvectors. Solution of systems of linear algebraic equa

**Solution of the normal equations**

⭐

The mapping of the sphere with diameter 1 via the matrix $\mathbf{A}^T\mathbf{A}$;
$$\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_{n-1} >> \lambda_n > 0$$



The mapping of the sphere with diameter 1 via the matrix $(\mathbf{A}^T\mathbf{A})^{-1}$;
$$\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_{n-1} >> \lambda_n > 0$$

Matrix equations, inverse of a matrix. Eigenvalues and eigenvectors of a matrix, generalized eigenvectors. Solution of systems of linear algebraic equa
○○○○○○○○○○○○○○○○

Solution of the normal equations

# Numerical solution of the normal equations

How large is in this case  the error of the computation   $E = \|\mathbf{A}^\mathrm{T}\mathbf{A}\mathbf{x} - \mathbf{A}^\mathrm{T}\mathbf{b}\|$ ?

In the finite computer arithmetic is, in general, the vector $\mathbf{A}^\mathrm{T}\mathbf{A}\mathbf{x}$ arbitrarily incorrect, because ones lost digits can't be gain back. It means that all components except the last one are damaged or even lost.

Numerical solution       $\mathbf{x} = (\mathbf{A}^\mathrm{T}\mathbf{A})^{-1}\mathbf{A}^\mathrm{T}\mathbf{b}$ :

- ★ we apply the Choleski decomposition to the symmetric positive definite matrix $\mathbf{A}^\mathrm{T}\mathbf{A}$.
  Disadvantage: it is necessary explicitly compute the matrix $\mathbf{A}^\mathrm{T}\mathbf{A}$, i.e., we have to compute many dot products that means that already the matrix $\mathbf{A}^\mathrm{T}\mathbf{A}$ can be computed with a quite large error. The advice is to apply a method that doesn't need the matrix $\mathbf{A}^\mathrm{T}\mathbf{A}$ to be explicitly given but works only with the matrix $\mathbf{A}$.
- via iterative methods.

The lost of the digits of the numerical computation is characterized by a so called condition number of the matrix. It is in this case equal to

$$\kappa(\mathbf{A}^\mathrm{T}\mathbf{A}) = \frac{\lambda_1}{\lambda_n}.$$

Matrix equations, inverse of a matrix. Eigenvalues and eigenvectors of a matrix, generalized eigenvectors. Solution of systems of linear algebraic equ
○○○○○○○○○○○○○○○○○
**Solution of the normal equations**

### Example

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

$$\begin{pmatrix} 10 & 7 & 8 & 7 \\ 7 & 5 & 6 & 5 \\ 8 & 6 & 10 & 9 \\ 7 & 5 & 9 & 10 \end{pmatrix} \begin{pmatrix} u_1 \\ u_2 \\ u_3 \\ u_4 \end{pmatrix} = \begin{pmatrix} 32 \\ 23 \\ 33 \\ 31 \end{pmatrix} \qquad \Rightarrow \quad \mathbf{x} = (1, 1, 1, 1)^\mathrm{T}$$

$$\mathbf{A}\widehat{\mathbf{x}} = \widehat{\mathbf{b}}$$

$$\begin{pmatrix} 10 & 7 & 8 & 7 \\ 7 & 5 & 6 & 5 \\ 8 & 6 & 10 & 9 \\ 7 & 5 & 9 & 10 \end{pmatrix} \begin{pmatrix} u_1 \\ u_2 \\ u_3 \\ u_4 \end{pmatrix} = \begin{pmatrix} 32,1 \\ 22,9 \\ 33,1 \\ 30,9 \end{pmatrix} \qquad \Rightarrow \quad \widehat{\mathbf{x}} = (9, 2; -12, 6; 4, 5; -1, 1)^\mathrm{T}$$

$$\widehat{\mathbf{A}}\widetilde{\mathbf{x}} = \mathbf{b}$$

$$\begin{pmatrix} 10 & 7 & 8,1 & 7,2 \\ 7,08 & 5,04 & 6 & 5 \\ 8 & 5,98 & 9,98 & 9 \\ 6,99 & 4,99 & 9 & 9,98 \end{pmatrix} \begin{pmatrix} u_1 \\ u_2 \\ u_3 \\ u_4 \end{pmatrix} = \begin{pmatrix} 32 \\ 23 \\ 33 \\ 31 \end{pmatrix} \qquad \Rightarrow \quad \widetilde{\mathbf{x}} = (-5, 79; 12, 02; -1.57, 2.57)^\mathrm{T}$$

The relative error:

$$\varepsilon_{rel}(\widehat{\mathbf{b}}) \quad = \quad \frac{\|\widehat{\mathbf{b}} - \mathbf{b}\|}{\|\mathbf{b}\|} = 0,003, \qquad \varepsilon_{rel}(\widehat{\mathbf{x}}) = 8,2$$

$$\varepsilon_{rel}(\widetilde{\mathbf{A}}) \quad = \quad \frac{\|\widetilde{\mathbf{A}} - \mathbf{A}\|}{\|\mathbf{A}\|} = 0,009, \qquad \varepsilon_{rel}(\widetilde{\mathbf{x}}) = 6,64$$

The matrix **A** is symmetric, $\det(\mathbf{A}) = 1$, but the condition number is $\kappa(\mathbf{A}) = 4488$

Matrix equations, inverse of a matrix. Eigenvalues and eigenvectors of a matrix, generalized eigenvectors. Solution of systems of linear algebraic equa
○○○○○○○○○○○○○○○

The solution of the normal equations by singular value decomposition

## ★ The normal equations and singular value decomposition

Let $\mathbf{A} = \mathbf{U}\mathbf{S}\mathbf{V}^{\mathrm{T}}$, $\quad \mathbf{x} = \arg \min_{\mathbf{y} \in \mathbb{R}^n} ||\mathbf{A}\mathbf{y} - \mathbf{b}|| \quad \Longrightarrow \quad \mathbf{S}\mathbf{V}^{\mathrm{T}}\mathbf{x} = \mathbf{U}^{\mathrm{T}}\mathbf{b}$, where
$\mathbf{S} = \mathrm{diag}(\sigma_1, \ldots, \sigma_p)$, $p = \min(m, n)$. Let the rank h($\mathbf{A}$)$= r < p$. Then
$\sigma_{r+1} = \cdots = \sigma_p = 0$, the matrix $S$ is singular and the inversion doesn't exist.
But if we multiply the second equation by the matrix $\mathbf{S}^+$ from the left

$$\mathbf{S}^+ = \begin{pmatrix} \mathbf{D}^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix}, \quad \mathbf{D}^{-1} = \begin{pmatrix} \dfrac{1}{\sigma_1} & & \\ & \ddots & \\ & & \dfrac{1}{\sigma_r} \end{pmatrix},$$

we obtain the system $\mathbf{V}^{\mathrm{T}}\mathbf{x} = \mathbf{S}^+\mathbf{U}^{\mathrm{T}}\mathbf{b}$ with the orthogonal matrix $\mathbf{V}^{\mathrm{T}}$. The
matrix $\mathbf{S}^+$ is so called Moore–Penrose pseudoinversion of the matrix $\mathbf{S}$.
Here, we will not study pseudoinversions.
A comparison:

$$\kappa(\mathbf{A}) = \frac{\sigma_1}{\sigma_r}, \quad \kappa(\mathbf{A}^{\mathrm{T}}\mathbf{A}) = \frac{\lambda_1}{\lambda_r} = \left(\frac{\sigma_1}{\sigma_r}\right)^2 = (\kappa(\mathbf{A}))^2 \quad \Longrightarrow$$

by the direct solution of the normal equations we lost two times more valid
digits then if we apply the singular value decomposition.

Matrix equations, inverse of a matrix. Eigenvalues and eigenvectors of a matrix, generalized eigenvectors. Solution of systems of linear algebraic equa
○○○○○○○○○○○○○○○○○

**The solution of the normal equations by singular value decomposition**

## Singular expansion of the matrix

The equation $\mathbf{A} = \mathbf{U}\mathbf{S}\mathbf{V}^{\mathrm{T}}$ can be rewritten as a sum of "singular terms" that correspond to the matrix $\mathbf{A}$:

$$\mathbf{A} = \sum_{i=1}^{r} \sigma_i \mathbf{u}_i \mathbf{v}_i^{\mathrm{T}}, \quad \mathrm{h}(\mathbf{A}) = r, \quad \mathrm{h}(\mathbf{u}_i \mathbf{v}_i^{\mathrm{T}}) = 1 \implies$$

$$\mathbf{x} \in \mathbb{R}^n \implies \mathbf{A}\mathbf{x} = \sum_{i=1}^{r} \sigma_i \mathbf{u}_i \mathbf{v}_i^{\mathrm{T}} \mathbf{x} = \sum_{i=1}^{r} (\mathbf{v}_i^{\mathrm{T}} \mathbf{x} \sigma_i) \mathbf{u}_i \ \ldots$$

a linear combination of the vectors $\mathbf{u}_i$, $i = 1, \ldots, r$.

### Application: data compression

Let the matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ contains the measured data. We are looking for an approximation of this matrix by the matrix $\mathbf{B}$ such that the rank of the matrix $\mathbf{B}$, $h(\mathbf{B}) = k < \min(m, n)$ and the matrix $\mathbf{B}$ contains all important information from the data in the matrix $\mathbf{A}$. For example if we would like to have the rank of the matrix $\mathbf{B}$ equal to 1, we just set $\mathbf{B} = \sigma_1 \mathbf{u}_1 \mathbf{v}_1^{\mathrm{T}}$.

Be careful! Different choice of $k$ will influence the quality of the results.

Matrix equations, inverse of a matrix. Eigenvalues and eigenvectors of a matrix, generalized eigenvectors. Solution of systems of linear algebraic equa
○○○○○○○○○○○○○○○○

Data compresion – an example

# ★ An example

**Example**    We want to digitalize a photo in such a way that we replace the photo by a matrix $24 \times 24$ pixels. The elements of the matrix are 0 (black box) or 1 (white box). We set the criterion for zero singular value to be smaller then $10^{-4}$ and obtain 16 nonzero singular values. All others are with this precision equal to 0:

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| $9,5403$ | $6,6288$ | $5,6369$ | $3,4756$ | $2,7385$ | $2,2023$ | $1,5835$ | $1,5566$ |
| $1,4207$ | $1,2006$ | $0,9905$ | $0,9258$ | $0,7479$ | $0,6744$ | $0,6122$ | $0,4698$ |

We are looking for such $k$ that the relative error will not be greater then 10. The relative error is
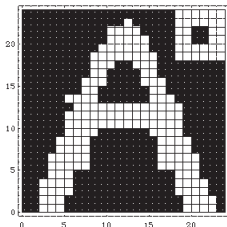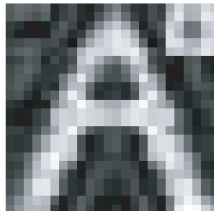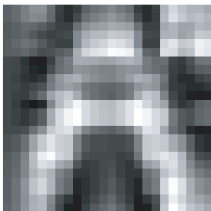
$$e(k) = 1 - \sqrt{\frac{\sum_{i=1}^{k} \sigma_i^2}{\sum_{i=1}^{16} \sigma_i^2}}$$

In particular, $e(2) = 0.18$, $e(3) = 0.09 \implies$ three "singular terms" of the matrix $\mathbf{A}$, $\implies$

$$\mathbf{B} = \sum_{i=1}^{3} \sigma_i \mathbf{u}_i \mathbf{v}_i^{\mathrm{T}}, \quad h(\mathbf{B}) = 3$$

Matrix equations, inverse of a matrix. Eigenvalues and eigenvectors of a matrix, generalized eigenvectors. Solution of systems of linear algebraic equa

○○○○○○○○○○○○○○○○○○○

**Data compresion – an example**



The original photo

Matrix equations, inverse of a matrix. Eigenvalues and eigenvectors of a matrix, generalized eigenvectors. Solution of systems of linear algebraic equa

**Data compresion – an example**

⭐



$k = 3$; $k = 5$; $k = 5$, elements of **B** are rounded to 0 or 1

## Recommended literature

- Ford W.: Numerical Linear Algebra with Applications Using MATLAB. Elsevier, 2015 .
- Golub G., Van Loan Ch: Matrix computations, The Johns Hopkins University press, Baltimore, 1996 (third edition) .
- Horn R. A. , Johnson Ch. R.: Matrix analysis, Cambridge University Press, 1985 .
- Kubíček M., Dubcová M., Janovská D.: Numerical methods and algorithms, http://old.vscht.cz/mat/Ang/NM-Ang/NM-Ang.pdf
- Póta G.: Mathematical Problems for Chemistry Students, Elsevier, 2006, ISBN 13:978-0-444-52793-6 (pbk.).
- Strang G.: Differential Equations and Linear Algebra, Department of Mathematics, Messsachusetts Institute of Technology, Wellesley-Cambridge Press, 2014.
- Trefethen N., Bau D.: Numerical linear algebra, SIAM Philadelphia, 1997 . ISBN 978-0-898713-61-9 (pbk.).
- Wilkinson J. H., Reinsch C.: Handbook for Automatie Computation. Vol. II Linear Algebra, Springer-Verlag, Berlin, Heidelberg, New York, 1971