# Numerical Calculus

## Course 3

## 2025

*Anca Ignat*

# Binary machine numbers[1]

In 1985 IEEE published a report entitled Binary Floating Point Arithmetic Standard 754-1985 and un update in 2008, IEEE 754-2008, that provides standards for binary and decimal floating point numbers, formats for data interchange, algorithms for rounding arithmetic operations, and handling exceptions. These standards are followed by computer manufacturers that use floating-point hardware.

[1] R.L.Burden, J.D. Faires, A.M. Burden, *Numerical Analysis*, 10th ed., Cengage Learning, Boston, USA

A 64-bit binary representations is done in the following way: the first bit is the sign bit, the next 11 bits represent the exponent $c$ which are followed by 52 bits that contain information regarding the binary fraction part, $f$, also called mantissa:

$$(-1)^s \, 2^{c-1023} (1+f) \, .$$

0 10000000011 1011100100010000000000000000000000000000000000000000 =
27.56640625

[27.56640624999999822364316059974953532218933105468750,
27.56640625000000177635683940025046467781066894531250).

The smallest positive number that can be represented using this type of representation, is when $s = 0, \ c = 1, f = 0$, that is:

$$z = 2^{-1022}(1+0) \approx 0.22251 \times 10^{-307}$$

and the biggest one is for $s = 0, c = 2046, f = 1 - 2^{-52}$

$$Z = 2^{1023}(2 - 2^{-52}) \approx 0.17977 \times 10^{309}.$$

The numbers that occur in computations and are smaller than $z$ are usually set to 0 (*underflow*) and those that are bigger than $Z$ usually stop the computations (*overflow*).

Note that number 0 has two representations: $s = 0, c = 1, f = 0$ and $s = 1, c = 1, f = 0$.

## Decimal Machine Numbers

Assume that the machine numbers are represented in the normalized decimal floating-point form:

$$\pm 0.d_1 d_2 \ldots d_k \times 10^n \quad 1 \leq d_1 \leq 9, \quad 0 \leq d_i \leq 9, i = 2, \ldots, k \quad -$$

This are $k$-digit decimal machine numbers. Any positive real number $y$ can be expressed by:

$$y = 0.d_1 d_2 \ldots d_k d_{k+1} d_{k+2} \ldots \times 10^n$$

A $k$-digit representation of $y$ can be obtained by a simple chopping operation:

$$fl(y) = 0.d_1 d_2 \ldots d_k \times 10^n \quad .$$

Another method, is the rounding operations:

$$fl(y) = 0.\delta_1 \delta_2 \ldots \delta_k \times 10^n$$

If $d_{k+1} \geq 5$ we add $1$ to $d_k$ to obtain $fl(y)$ (*round up*), otherwise

a $k$-digit chop is performed (*round down*).

A number $r^*$ approximates the number $r$ to $t$ significant digits if $t$ is the largest nonnegative integer for which:

$$\frac{\left|r - r^*\right|}{|r|} \leq 5 \times 10^{-t} \ .$$

For chopping we have:

$$\left|\frac{y - fl(y)}{y}\right| \leq 10^{-k+1}$$

And for rounding:

$$\left|\frac{y - fl(y)}{y}\right| \leq 0.5 \times 10^{-k+1}.$$

# Finite-Digit Arithmetic

$$x +_c y = fl(fl(x) + fl(y))$$

$$x -_c y = fl(fl(x) - fl(y))$$

$$x \times_c y = fl(fl(x) \times fl(y))$$

$$x \div_c y = fl(fl(x) \div fl(y))$$

# Sources of error in numerical computations

1. Errors in <u>input data</u>:

    - systematic errors or temporary perturbations in measurements,

    - round-off errors whenever dealing with irrational numbers, or with infinite decimal representation: 1/3, $\pi$, 1/7,…

2. <u>Round-off</u> errors during computations:

    - the result of elementary operations is not, always, exact

3. <u>Discretization</u> errors:

- the solution of the solved problem is the limit of a sequence, the sum of an infinite series,

- nonlinear functions are approximated by linear functions, computing the approximate derivative or integral for a function

4. <u>Simplifications</u> in the mathematical model

- idealizations, ignoring some parameters

5. <u>Human</u> errors, bugs in numerical software.

# Absolute Error, Relative Error

$a$ – exact value

$\tilde{a}$ – approximate value

*Absolute error* : $a - \tilde{a}$ or $|a - \tilde{a}|$ or $\|a - \tilde{a}\|$

$$a = \tilde{a} \pm \Delta_a , \ |a - \tilde{a}| \leq \Delta_a$$

*Relative error*: $a \neq 0$ $\quad \dfrac{a - \tilde{a}}{a}$ or $\dfrac{|a - \tilde{a}|}{|a|}$ or $\dfrac{\|a - \tilde{a}\|}{\|a\|}$

$$\dfrac{|a - \tilde{a}|}{|a|} \leq \delta_a \quad (\delta_a \text{ usually expressed in } \%).$$

If we have the following approximations: 1kg ±5g, 50g±5g, the absolute errors are the same, but for the first quantity the relative error is 0,5% and for the second, the relative error is **10%**.

$$a_1 = \tilde{a}_1 \pm \Delta_{a_1}, a_2 = \tilde{a}_2 \pm \Delta_{a_2},$$

$$a_1 \pm a_2 = (\tilde{a}_1 \pm \tilde{a}_2) \pm \left(\Delta_{a_1} \pm \Delta_{a_2}\right)$$

$$\Delta_{a_1 + a_2} \leq \Delta_{a_1} + \Delta_{a_2}.$$

$a_1$ has $\delta_{a_1}$ relative error and $a_2$ has relative error $\delta_{a_2}$ :

$$a = a_1 * a_2 \text{ or } \frac{a_1}{a_2} \rightarrow \delta_a = \delta_{a_1} + \delta_{a_2}.$$

## Conditioning ⟵⟶ stability

The conditioning of a problem describes the sensibility/variation of the solution with respect to the input data, assuming that all the computations are exact (regardless of the algorithm employed for solving the problem).

Let $x$ be the exact input data, $\tilde{x}$ an approximation for the input data, $P(x)$ the exact solution of the problem, and $P(\tilde{x})$ the solution with $\tilde{x}$ as input. We assume that the

computations were exact in obtaining the solutions *P(x)* and *P($\tilde{x}$)* .

O problem is considered ***ill-conditioned*** if the difference between *P(x)* and *P($\tilde{x}$)* is large even if the relative error in input data $\dfrac{\| x - \tilde{x} \|}{\| x \|}$ is small.

The numerical conditioning of a problem is expressed as the ratio between output and input relative errors:

$$k(x) = \frac{\dfrac{\parallel P(x) - P(\tilde{x}) \parallel}{\parallel P(x) \parallel}}{\dfrac{\parallel x - \tilde{x} \parallel}{\parallel x \parallel}} \quad \text{for } x \neq 0 \text{ and } P(x) \neq 0$$

A small value for $k(x)$ characterizes a well-conditioned problem. Conditioning is a local property (it depends on the input $x$). A problem is well-conditioned if it is well-conditioned for all input data.

Relative error in the output $\approx$

Condition number $\times$ Relative error in the input

Consider Wilkinson polynomial:

$$w(x) = (x-1)(x-2)\cdots(x-20) = x^{20} - 210x^{19} + P_{18}(x)$$

If one changes the coefficient (**-210**) of $x^{19}$ with:

$$-210 - 2^{-23} = -210.0000001192$$

the new polynomial's roots (with 5 significant digits) are:

**1.00000 2.00000 3.00000 4.00000 5.00000 6.00001 6.99970 8.00727**

**8.91725 20.84691 10.09527 $\pm$ 0.64350$i$ 11.79363 $\pm$ 1.65233$i$**

**13.99236 $\pm$ 2.51883$i$ 16.73074 $\pm$ 2.81262$i$ 19.50244 $\pm$ 1.94033$i$**

For solving a problem $P$, an algorithm $\tilde{P}$ is implemented and run on a computer. Taking into account the inexact storing of data and inaccurate computations the exact solution and the computed one can be different:

$$P(x) \neq \tilde{P}(x)$$

*Numerical stability* measures the magnitude of the numerical errors after running the implemented algorithm, assuming that the input data are exact, i.e., $\| P(x) - \tilde{P}(x) \|$ or

$$\frac{\| P(x) - \tilde{P}(x) \|}{\| P(x) \|}.$$

A *numerically stable algorithm* is an algorithm for which the relative error is the same order as the rounding error.

A ***numerically stable algorithm*** used for solving a **well-conditioned problem** yields results with very good precision.

An algorithm $\tilde{P}$ for solving problem $P$ is numerically stable if one of the following condition is fulfilled:

1. $\tilde{P}(x) \approx P(x)$ for all input $x$;


2. there exists $\tilde{x}$ near $x$, s.t. $\tilde{P}(x) \approx P(\tilde{x})$

$$x = \text{exact data,}$$

$$P(x) = \text{exact solution using exact input,}$$

$$\tilde{P}(x) = \text{„}computed\text{''} \text{ solution using algorithm } \tilde{P} \text{ with exact}$$

$$\text{input data}$$

# Solving linear systems of equations

- 1900 B.C., Babylon – first problems related to simultaneous linear equations

- 300 B.C. Babylon – problem on clay tablets:

*" There are two fields whose total area is 1800 square yards. One produces grain at the rate of 2/3 of a bushel per square yard while the other produces grain at the rate of 1/2 a bushel per square yard. If the total yield is 1100 bushels, what is the size of each field?"*

- 200-100 B.C. China – *Nine Chapters on the Mathematical Art* – a method very similar to Gaussian Elimination for solving a 3-dimensional linear system of equations (*„There are three types of corn, of which three bundles of the first, two of the second, and one of the third make 39 measures. Two of the first, three of the second and one of the third make 34 measures. And one of the first, two of the second and three of the third make 26 measures. How many measures of corn are contained of one bundle of each type?"*)

- 1545, Cardan – in *Ars Magna*, proposes a rule (*regula de modo)* for solving a system of 2 equations with 2 unknowns (similar to Cramer's rule)

- 1683, Seki Kowa, Japan – introduces the idea of „*determinant*"- „*Method of solving the dissimulated problems*". He computes measures that we today call determinants for matrices **2x2**, **3x3**, **4x4**, **5x5** related to solving equations but not systems of linear equations.

- 1683, Leibniz in a letter to l'Hôpital explains that the system of linear equations:

$$10+11x+12\,y=0$$
$$20+21x+22\,y=0$$
$$30+31x+32\,y=0$$

has a solution because:

**10\*21\*32+11\*22\*30+12\*20\*31=10\*22\*31+11\*20\*32+12\*21\*30**

(the condition that the determinant of the coefficients matrix must be 0).

Leibniz was convinced that good mathematical notation was the key to progress so he experimented with different notation for coefficient systems. His unpublished manuscripts contain more than 50 different ways of writing coefficient systems which he worked on during a period of 50 years beginning in 1678. Leibniz uses the notion of „*resultant*" instead of determinant and he proved Cramer's rule for „resultants". He knew that every determinant could be expanded using any column (Laplace expansion).

- 1750, Cramer introduces a formula based on determinants for solving a system of linear equations – *Cramer's rule* – *„Introduction in the analysis of algebraic curves"* (he presents an algorithm for general *n x n* systems)*:

  *,One finds the value of each unknown by forming **n** fractions of which the common denominator has as many terms as there are permutations of **n** things'*

- 1764 Bezout, 1771 Vandermonde, 1772 Laplace – rules for computing determinants

- 1773 Lagrange – first implicit use of matrices related to bilinear forms that appear in the problem of optimization for real functions depending on 2 or more variables (he wanted to characterize the optimum points for functions depending on 2 or more variables)

- 1800-1801, Gauss introduces for the first time the term *„determinant"* (determines the properties of the quadratic form) – *Disquisitiones arithmeticae(1801)*; he describes the matrix multiplication and matrix inverse notions in the context of quadratic forms and the associated bi-dimensional table of coefficients. Gauss proposed the *Gaussian elimination* method for solving a linear system with 6 equations in six unknowns, when he studied the orbit of Pallas asteroid.

- 1812, Cauchy uses the term of „*determinant*" in the modern sense.

- 1826, Cauchy finds the eigenvalues and deduces results for diagonalizing a matrix. He introduces the idea of similar matrices and he proves that these matrices have the same characteristic equations. He also proves that every symmetric real matrix is diagonalizable.

- 1850, Sylvester introduces for the first time the term *matrix* (from Latin, „uterus, womb" – a place something is formed or produced, „*an oblong arrangement of terms*")
- 1855, Cayley – matrix algebra, first abstract definition of a matrix. He studies linear transformations and their composition which leads him to the well-known matrix operations (addition, multiplication, scalar multiplication, inverse)

- 1858, Cayley in *Memoir on the theory of matrices:* *„there would be many things to say about this theory of matrices which should, it seems to me, precede the theory of determinants"*

- Jordan (1870 – Treatise on substitutions and algebraic equations – Jordan canonical form), Frobenius (1878 – On linear substitutions and bilinear forms, rank of a matrix)

- 1890, Weierstrass – On determinant theory, the axiomatic definition for the determinant

- 1925, Heisenberg reinvents matrix algebra for quantum mechanics

- 1947, vonNeuman & Goldstine introduce conditioning numbers when analysing rounding errors

- 1948, Turing introduces the $LU$ decomposition of a matrix i

- 1958, Wilkinson - $QR$ factorization …

**Proposition**

Let $A \in \mathbb{R}^{n \times n}$ for which there exists a matrix natural norm s.t. $\|A\| < 1$. Then there exist the matrices $(I_n \pm A)^{-1}$, and the following relations hold:

$$\frac{1}{1+\|A\|} \le \left\|(I \pm A)^{-1}\right\| \le \frac{1}{1-\|A\|}.$$

# Error evaluation in solving linear systems

(linear systems of equations conditioning)

Let $A \in \mathbb{R}^{n \times n}$, $b \in \mathbb{R}^{n}$, $x \in \mathbb{R}^{n}$ defining the linear system of equations:

$$Ax = b$$

$A$ **non-singular** $\Leftrightarrow$ $\det A \neq 0$ $\Rightarrow$ $\exists$ **solution** $x = A^{-1}b$

We consider the following errors in data input:

- $\Delta A \in \mathbb{R}^{n \times n}$ **absolute error for $A$**;

- $\Delta b \in \mathbb{R}^{n}$ **absolute error for $b$**;

Instead of solving the linear system $Ax=b$, the following system is solved:

$$(A + \Delta A)\tilde{x} = b + \Delta b$$

one obtains the solution $\tilde{x}$:

$$\tilde{x} = x + \Delta x$$

The following problems must be addressed:

1.   $A$ non-singular, $\Delta A = ?$ s.t. $A + \Delta A$ is also non-singular?

2.   Assume $A$ and $A + \Delta A$ non-singular, which are the relations between the relative errors $\dfrac{\|\Delta A\|}{\|A\|}, \dfrac{\|\Delta b\|}{\|b\|}$ and $\dfrac{\|\Delta x\|}{\|x\|}$?

1. Assume $A$ non-singular:

$$A + \Delta A = A\left(I_n + A^{-1}\Delta A\right) \rightarrow$$

$A + \Delta A$ **non-singular** $\Leftrightarrow \left(I_n + A^{-1}\Delta A\right)$ **non-singular**

*Proposition 1*

Let $A$ be non-singular and $\|\Delta A\| < \dfrac{1}{\|A^{-1}\|}$. Then $I + A^{-1}\Delta A$ is

non-singular and the following relation holds:

$$\left\|\left(I_n + A^{-1}\Delta A\right)^{-1}\right\| \leq \dfrac{1}{1 - \|A^{-1}\| \cdot \|\Delta A\|}.$$

Proof. We have:

$$\left\| \Delta A \right\| < \frac{1}{\left\| A^{-1} \right\|} \Rightarrow \left\| A^{-1} \Delta A \right\| \leq \left\| A^{-1} \right\| \cdot \left\| \Delta A \right\| < 1 \overset{\textbf{Prop}}{\Rightarrow} \exists \left( I + A^{-1} \Delta A \right)^{-1}$$

$$\left\| \left( I + A^{-1} \Delta A \right)^{-1} \right\| \leq \frac{1}{1 - \left\| A^{-1} \Delta A \right\|} \leq \frac{1}{1 - \left\| A^{-1} \right\| \cdot \left\| \Delta A \right\|}.$$

Assume $A$ non-singular and $\left\| \Delta A \right\| < \dfrac{1}{\left\| A^{-1} \right\|}$.

$$\frac{\left\| \Delta x \right\|}{\left\| x \right\|} \leq \frac{\left\| A^{-1} \right\| \left\| A \right\|}{1 - \left\| A^{-1} \right\| \left\| \Delta A \right\|} \left( \frac{\left\| \Delta b \right\|}{\left\| b \right\|} + \frac{\left\| \Delta A \right\|}{\left\| A \right\|} \right).$$

36

$$\left(A+\Delta A\right)\left(x+\Delta x\right)=b+\Delta b \Rightarrow \left(A+\Delta A\right)\Delta x+Ax+\left(\Delta A\right)x=b+\Delta b \Rightarrow$$

$$A\left(I+A^{-1}\Delta A\right)\Delta x=\Delta b-\left(\Delta A\right)x \Rightarrow \Delta x=\left(I+A^{-1}\Delta A\right)^{-1}A^{-1}\left[\Delta b-\left(\Delta A\right)x\right]\Rightarrow$$

$$\left\|\Delta x\right\|\le\left\|\left(I+A^{-1}\Delta A\right)^{-1}\right\|\left\|A^{-1}\right\|\left(\left\|\Delta b\right\|+\left\|\Delta A\right\|\left\|x\right\|\right)\Rightarrow$$

$$\frac{\left\|\Delta x\right\|}{\left\|x\right\|}\le\frac{\left\|A^{-1}\right\|}{1-\left\|A^{-1}\right\|\left\|\Delta A\right\|}\left(\frac{\left\|\Delta b\right\|}{\left\|x\right\|}+\left\|\Delta A\right\|\right) \qquad (1)$$

From $Ax=b$ we get $\left\|b\right\|\le\left\|A\right\|\left\|x\right\|\Rightarrow\dfrac{1}{\left\|x\right\|}\le\dfrac{\left\|A\right\|}{\left\|b\right\|}$ , using this

relation in (1) we have:

$$\frac{\left\|\Delta x\right\|}{\left\|x\right\|}\le\frac{\left\|A^{-1}\right\|\left\|A\right\|}{1-\left\|A^{-1}\right\|\left\|\Delta A\right\|}\left(\frac{\left\|\Delta b\right\|}{\left\|b\right\|}+\frac{\left\|\Delta A\right\|}{\left\|A\right\|}\right).$$

*k(A) = ||A⁻¹|| ||A|| **conditioning number*** for matrix *A*.

**Proposition 2**

If matrix $A$ is non-singular and $\|\Delta A\| < \dfrac{1}{\|A^{-1}\|}$ then:

$$\frac{\|\Delta x\|}{\|x\|} \le \frac{k(A)}{1 - k(A) \cdot \dfrac{\|\Delta A\|}{\|A\|}} \left( \frac{\|\Delta b\|}{\|b\|} + \frac{\|\Delta A\|}{\|A\|} \right).$$

From $I_n = A\,A^{-1}$ one gets $1 = \|I_n\| \le \|A\|\|A^{-1}\| = k(A).$

*k(A) ≥ 1,* ∀*A* but it depends on the natural matrix norm.

A matrix $A$ for which the conditioning number is large is called **ill-conditioned (k(A) ,large').**

$Ax=b$ **with** $k(A)$ **large** $\rightarrow$ $\dfrac{\|\Delta x\|}{\|x\|}$ the relative error in the

computed solution can be very large even if the relative errors

$\dfrac{\|\Delta b\|}{\|b\|}$ **and** $\dfrac{\|\Delta A\|}{\|A\|}$ are small.

Let $A$ be a symmetric, non-singular matrix $A = A^T$, $\det A \neq 0$.

Using the spectral norm (generated by the Euclidean vector norm) we have:

$$\|A\|_2 = \sqrt{\rho\left(A^T A\right)} = \sqrt{\rho\left(A^2\right)}$$

$$k\left(A\right) = \|A\|_2 \cdot \|A^{-1}\|_2$$

A symmetric matrix $A$ has only real eigenvalues $\lambda_1, \lambda_2, ..., \lambda_n$,

$A^2$ has eigenvalues $\lambda_1^2, \lambda_2^2, ..., \lambda_n^2$, $A^{-1}$ has eigenvalues

$$\frac{1}{\lambda_1}, \frac{1}{\lambda_2}, ....., \frac{1}{\lambda_n}.$$

$$|\lambda_1| \le |\lambda_2| \le .... \le |\lambda_n| \Rightarrow \rho(A) = |\lambda_n| \text{ and } \rho(A^{-1}) = \frac{1}{|\lambda_1|}$$

$$A = A^T \rightarrow \|A\|_2 = \rho(A) = |\lambda_n|, \|A^{-1}\|_2 = \rho(A^{-1}) = \frac{1}{|\lambda_1|},$$

$$k_2(A) = \|A\|_2 \cdot \|A^{-1}\|_2 = \frac{|\lambda_n|}{|\lambda_1|} \text{ spectral conditioning number.}$$

$A$ orthogonal matrix $\rightarrow k_2(A) = 1$

$$A^T A = A \cdot A^T = I_n \Rightarrow A^{-1} = A^T$$

$$\|A\|_2 = \sqrt{\rho(A^T A)} = \sqrt{\rho(I)} = 1 = \|A^T\|_2$$

$$k_2(A) = \|A\|_2 \cdot \|A^{-1}\|_2 = \|A\|_2 \cdot \|A^T\|_2 = 1,$$

Almost singular matrix with good conditioning number

$$A = \mathbf{diag}\,[1, 0.1, 0.1, ..., 0.1] \in \mathbb{R}^{100 \times 100} \Rightarrow \det A = 1 \cdot (0.1)^{99} = 10^{-99}$$

$$\| A \|_2 = 1 \;,\; \| A^{-1} \|_2 = 10 \Rightarrow k_2(A) = \| A \|_2 \; \| A^{-1} \|_2 = 10$$

Very ill-conditioned matrix with non-zero determinant
(**det** $A$=**1**)

$$A = \begin{pmatrix} 1 & 2 & 0 & \cdots & 0 \\ 0 & 1 & 2 & \cdots & 0 \\ \vdots & & & & \\ 0 & 0 & 0 & \cdots & 2 \\ 0 & 0 & 0 & \cdots & 1 \end{pmatrix},$$

$$A^{-1} = \begin{pmatrix} 1 & -2 & 4 & \cdots & (-2)^{i-1} & \cdots & (-2)^{n-1} \\ 0 & 1 & -2 & \cdots & (-2)^{i-2} & \cdots & (-2)^{n-2} \\ \vdots & & & & & & \\ 0 & 0 & 0 & \cdots & 0 & \cdots & 1 \end{pmatrix}$$

$$\| A \|_\infty = \| A \|_1 = 3 \quad ,$$

$$\| A^{-1} \|_\infty = \| A^{-1} \|_1 = 1 + 2 + \cdots + 2^{n-1} = 2^n - 1$$

$$n = 100 \;\Rightarrow\; k(A) = \| A \|_\infty \| A^{-1} \|_\infty = \| A \|_1 \| A^{-1} \|_1 = 3 \cdot (2^{100} - 1)$$

$$\det A = 1$$

$$\begin{cases} x + y = 2 \\ x + 1.001y = 2 \end{cases} \quad, \quad \begin{cases} x + y = 2 \\ x + 1.001y = 2.001 \end{cases} \quad k(A) = 4002$$

$$x = 2, \ y = 0 \qquad\qquad x = 1, \ y = 1$$

$$\begin{cases} 400x - 201y = 200 \\ -800x + 401y = -200 \end{cases} \quad, \quad \begin{cases} 401x - 201y = 200 \\ -800x + 401y = -200 \end{cases}$$

$$x = -100, \ y = -200 \qquad\qquad x = 40000, \ y = 79800$$

$$k_2(A) = 2503 \qquad\qquad k_2(A) = 1002000$$

$$\begin{cases} 1.2969x + 0.8648\,y = 0.8642 \\ 0.2161x + 0.1441\,y = 0.1440 \end{cases} \qquad x = 2\,,\ y = -2 \qquad k_2(A) = 249730000$$

$$\bar{x} = 0.9911,\ \bar{y} = -0.4870\ ,$$

$$r = b - Az = \begin{pmatrix} 0.8642 \\ 0.1440 \end{pmatrix} - \begin{pmatrix} 1.2969 & 0.8648 \\ 0.1441 & 0.1441 \end{pmatrix} \begin{pmatrix} 0.9911 \\ -0.4870 \end{pmatrix} = \begin{pmatrix} 10^{-8} \\ -10^{-8} \end{pmatrix}$$

## Hilbert matrix

$$H = (h_{ij})_{i,j=1}^{n} \quad , \quad h_{ij} = \frac{1}{i+j-1} = \int_{0}^{1} x^{i+j-2} \, dx$$

$$k_2(H_n) \approx \frac{(\sqrt{2}+1)^{4(n+1)}}{2^{\frac{15}{4}} \sqrt{\pi n}} \sim e^{3.5n}$$

| n | $k_2(H_n)$ | n | $k_2(H_n)$ |
|---|---|---|---|
| 1 | 1 | 7 | $4.753 \cdot 10^8$ |
| 2 | 19.281 | 8 | $1.526 \cdot 10^{10}$ |
| 3 | $5.241 \cdot 10^2$ | 9 | $4.932 \cdot 10^{11}$ |
| 4 | $1.551 \cdot 10^4$ | 10 | $1.602 \cdot 10^{13}$ |
| 5 | $4.766 \cdot 10^5$ | 11 | $5.220 \cdot 10^{14}$ |
| 6 | $1.495 \cdot 10^7$ | 12 | $1.678 \cdot 10^{16}$ |

$$H^{-1} = (g_{ij}) \quad g_{ij} = \frac{(-1)^{i+j}}{(i+j-1)} \frac{(n+i-1)! \, (n+j-1)!}{\left[(i-1)!(j-1)!\right]^2 (n-i)!(n-j)!}$$