

A CLASS OF ARBITRARILY ILL CONDITIONED FLOATING-POINT MATRICES*

SIEGFRIED M. RUMP†

Abstract. Let \mathbb{F} be a floating-point number system with basis $\beta \geq 2$ and an exponent range consisting of at least the exponents 1 and 2. A class of arbitrarily ill conditioned matrices is described, the coefficients of which are elements of \mathbb{F} . Due to the very rapidly increasing sensitivity of those matrices, they might be regarded as "almost" ill posed problems.

The condition of those matrices and their sensitivity with respect to inversion is given by means of a closed formula. The condition is rapidly increasing with the dimension. For example, in the double precision of the IEEE 754 floating-point standard (base 2, 53 bits in the mantissa including implicit 1), matrices with $2n$ rows and columns are given with a condition number of approximately $4 \cdot 10^{32n}$.

Key words. condition number, sensitivity, ill conditioned, linear systems, floating-point number systems

AMS(MOS) subject classifications. 15A12, 65F05, 65G05

0. Introduction. It is a trivial fact that there are arbitrarily ill conditioned *real* matrices. In this paper we concentrate on matrices that are exactly representable in some floating-point number system \mathbb{F} . There is no restriction to the basis and only a trivial technical assumption on the exponent range of \mathbb{F} . For fixed \mathbb{F} there are finitely many square matrices with n rows and a maximum condition number less than ∞ for given n .

The well-known schemes for constructing ill-conditioned matrices suffer from the fact that for given \mathbb{F} only a few matrices are exactly representable in \mathbb{F} , say up to n_{\max} rows. For $n > n_{\max}$ rows the entries are getting "too big." For example, let

$$(Z_n)_{ij} := \frac{\binom{n+i-1}{i-1} \cdot n \cdot \binom{n-1}{n-j}}{i+j-1},$$

as proposed by Zielke. For single precision in the IEEE 754 floating-point format (base 2 with 24 bit in the mantissa including implicit 1), we have (using infinity norm)

$$n_{\max}(Z_n) = 10 \quad \text{with} \quad \|Z_{10}\| \cdot \|Z_{10}^{-1}\| \approx 2 \cdot 10^{14}.$$

From Pascal's triangle we get

$$(P_n)_{ij} := \binom{i+j-1}{i-1}$$

with

$$n_{\max}(P_n) = 15 \quad \text{with} \quad \|P_{15}\| \cdot \|P_{15}^{-1}\| \approx 1 \cdot 10^{16}.$$

The classical example for ill-conditioned matrices is Hilbert matrices, the ij th component of which is $1/(i+j-1)$. In order to make them exactly representable in a binary floating-point format, we may use their inverses, or we may multiply the entire matrix by $lcm(1, 2, \dots, 2n-1)$. We call the latter matrix H_n^* . Then

$$n_{\max}(H_n^{-1}) = 7 \quad \text{with} \quad \|H_7\| \cdot \|H_7^{-1}\| \approx 5 \cdot 10^8$$

* Received by the editors August 31, 1989; accepted for publication (in revised form) March 13, 1990.

† Informatik III, Technische Universität, 2100 Hamburg 90, Federal Republic of Germany.

consists only of components that are exactly representable in \mathbb{F} . Since (1.4) has infinitely many solutions, the class of matrices C_n defined by (1.6) consists of elements with an arbitrarily large number of rows.

2. Properties of the matrices. In this section some properties of the matrices defined by (1.6) will be studied. Here, no restrictions on k or σ with respect to β are necessary; our only assumptions are (1.5) and (1.4). In the following, especially, the assumption $0 \leq p_i, q_i < \sigma$ for $i = 0 \cdots n$ is not necessary.

Throughout this paper we use componentwise ordering of matrices, i.e., $A \leq B : \langle = \rangle a_{ij} \leq b_{ij}$ and the componentwise absolute value $|A| = (|A_{ij}|)$, which is again a matrix.

The condition number $\|C_n\| \cdot \|C_n^{-1}\|$ for the ∞ -norm will be calculated along with the sensitivity of C_n . Rohn, in [3], gave a nice definition of the sensitivity of a matrix C with respect to inversion: Let B be a matrix of relative distance less than or equal to α to C , i.e., $|B - C| \leq \alpha \cdot |C|$, then

$$s_{ij}^\alpha(C) := \max \left\{ \frac{|B_{ij}^{-1} - C_{ij}^{-1}|}{|C_{ij}^{-1}|}; |B - C| \leq \alpha \cdot |C| \right\},$$

provided $C_{ij}^{-1} \neq 0$ and

$$s_{ij}(C) := \lim_{\alpha \rightarrow 0^+} \frac{s_{ij}^\alpha(C)}{\alpha}.$$

In [3], Rohn proves an explicit formula for the sensitivity matrix $S = (s_{ij}(C))$:

$$(2.1) \quad s_{ij}(C) = \frac{(|C^{-1}| \cdot |C| \cdot |C^{-1}|)_{ij}}{|C^{-1}|_{ij}} \quad \text{for } C_{ij}^{-1} \neq 0.$$

LEMMA 1. $\det(C_0) = 1$, $\|C_0\|_\infty \|C_0^{-1}\|_\infty = (P + kQ)^2$, and $s_{ij}(C_0) = 4P^2 - 3$ for $i = j$ and $s_{ij}(C_0) = 4P^2 - 1$ for $i \neq j$.

Proof. For $n = 0$, (1.6) writes

$$C_0 = \begin{pmatrix} P & kQ \\ Q & P \end{pmatrix} \quad \text{with } C_0^{-1} = \begin{pmatrix} P & -kQ \\ -Q & P \end{pmatrix},$$

as follows from (1.4). Then the first two statements are obvious; for the third, a short computation yields

$$(s_{ij}(C_0)) = \begin{pmatrix} \zeta & \eta \\ \eta & \zeta \end{pmatrix} \quad \text{with } \zeta = P^2 + 3kQ^2, \quad \eta = 3P^2 + kQ^2. \quad \square$$

In the following we will show that for $n > 0$ the condition and sensitivity of C_n increase compared to those of C_0 .

For the rest of the paper we frequently use

$$(2.2) \quad C := C_n \in \mathbb{R}^{(2n+2) \times (2n+2)} \quad \text{with components } c_{ij}, 0 \leq i, j \leq 2n+1.$$

The indices of matrices start with 0 with the exception of A and B , to be defined later. Those are $(n + 1) \times n$ -matrices with row indices starting with σ and column indices starting with 1.

LEMMA 2. *The matrices C_n are not singular: $\det(C_n) = (-1)^n$.*

Proof. Define

$$(2.3) \quad s := (\sigma^n, \sigma^{n-1}, \dots, \sigma, 1)^t \in \mathbb{R}^{n+1}$$

and

$$(2.4) \quad x := \begin{pmatrix} \frac{P \cdot s}{-Q \cdot s} \end{pmatrix} = \begin{pmatrix} P \cdot \sigma^n \\ \vdots \\ P \cdot 1 \\ -Q \cdot \sigma^n \\ \vdots \\ -Q \cdot 1 \end{pmatrix} \in \mathbb{R}^{2n+2}.$$

Then

$$(2.5) \quad (p_n, \dots, p_0) \cdot s = P \quad \text{and} \quad (q_n, \dots, q_0) \cdot s = Q,$$

and using (2.2),

$$\sum_{v=0}^{2n+1} c_{0v} \cdot x_v = P^2 - kQ^2 = 1,$$

$$\sum_{v=0}^{2n+1} c_{1v} \cdot x_v = PQ - QP = 0 = \sum_{v=0}^{2n+1} c_{iv} \cdot x_v \quad \text{for } i \geq 2.$$

This means that x is the first column of C^{-1} and, especially,

$$(2.6) \quad (C^{-1})_{2n+1,0} = -Q.$$

Therefore $-Q = -\det(\bar{C})/\det(C)$ with

$$\bar{C} := \begin{pmatrix} q_n \cdots q_0 & p_n \cdots p_1 \\ \Sigma & 0 \\ 0 & \Sigma^* \end{pmatrix},$$

and

$$\Sigma := \begin{pmatrix} 1 & -\sigma & & & \\ & 1 & -\sigma & & \\ & & \dots & & \\ & & & 1 & -\sigma \end{pmatrix}, \quad \Sigma^* := \begin{pmatrix} 1 & -\sigma & & & \\ & 1 & -\sigma & & \\ & & \dots & & \\ & & & \dots & \\ & & & & 1 \end{pmatrix}.$$

But $\det(\bar{C}) = \det(\bar{\bar{C}})$ with

$$\bar{\bar{C}} := \begin{pmatrix} q_n \cdots q_0 \\ \Sigma \end{pmatrix}$$

and $\bar{\bar{C}} \cdot s = Q \cdot e$ with $e = (1, 0, \dots, 0)^t$. This implies that

$$(\bar{\bar{C}}^{-1})_{00} = \sigma^n / Q = \det(\hat{C}) / \det(\bar{\bar{C}})$$

with

$$\hat{C} := \begin{pmatrix} -\sigma & & & & \\ 1 & -\sigma & & & \\ & & \dots & & \\ & & & 1 & -\sigma \end{pmatrix} \in \mathbb{R}^{n \times n}, \quad \det(\hat{C}) = -(1)^n \cdot \sigma^n.$$

Therefore

$$\det(C) = \frac{\det(\bar{C})}{Q} = \frac{\det(\bar{\bar{C}})}{Q} = \frac{\det(\hat{C}) \cdot Q}{\sigma^n \cdot Q} = (-1)^n.$$

□

Next we calculate the inverse of $C = C_n$ explicitly. The first column is already given by (2.4), the second is given by

$$(2.7) \quad y := \begin{pmatrix} -k \cdot Q \cdot s \\ P \cdot s \end{pmatrix} \in \mathbb{R}^{2n+2}, \quad C \cdot y = (0, 1, 0, \dots, 0)^t.$$

Formulas (2.4) and (2.7) imply, especially, that $-Q$ and P are the first two elements of the last row of C^{-1} . Let

$$(2.8) \quad (-QP\alpha_n \cdots \alpha_1 \beta_n \cdots \beta_1) \in \mathbb{R}^{2n+2}$$

be the last row of C^{-1} . Then multiplication with the first $n + 1$ columns of C yields

$$(2.9) \quad \begin{aligned} -Q \cdot p_n + P \cdot q_n + \alpha_n &= 0, \\ -Q \cdot p_{n-1} + P \cdot q_{n-1} - \sigma \cdot \alpha_n + \alpha_{n-1} &= 0, \\ &\dots \\ -Q \cdot p_1 + P \cdot q_1 - \sigma \cdot \alpha_2 + \alpha_1 &= 0, \\ -Q \cdot p_0 + P \cdot q_0 - \sigma \cdot \alpha_1 &= 0. \end{aligned}$$

Setting $\alpha_0 = \alpha_{n+1} = 0$ by definition gives

$$(2.10) \quad -Q \cdot p_i + P \cdot q_i - \sigma \cdot \alpha_{i+1} + \alpha_i = 0 \quad \text{for } i = 0 \cdots n$$

and by successively adding the equations in (2.9), multiplied by σ , yields

$$(2.11) \quad \alpha_i = Q \cdot \sum_{v=i}^n p_v \cdot \sigma^{v-i} - P \cdot \sum_{v=i}^n q_v \cdot \sigma^{v-i} \quad \text{for } i = 1 \cdots n.$$

By treating the last $n + 1$ columns of C in the same way, we obtain

$$(2.12) \quad \begin{aligned} -k \cdot Q \cdot q_i + P \cdot p_i - \sigma \cdot \beta_{i+1} + \beta_i &= 0 \quad \text{for } i = 1 \cdots n, \\ -k \cdot Q \cdot q_0 + P \cdot p_0 - \sigma \cdot \beta_1 &= 1, \end{aligned}$$

setting $\beta_0 = \beta_{n+1} = 0$ by definition, and

$$(2.13) \quad \beta_i = P \cdot \sum_{v=i}^n p_v \cdot \sigma^{v-i} - k \cdot Q \cdot \sum_{v=i}^n q_v \cdot \sigma^{v-i} \quad \text{for } i = 1 \cdots n.$$

According to our assumption (1.5), $p_n \neq 0$ or $q_n \neq 0$ and

$$\sum_{v=i}^n p_v \cdot \sigma^{v-i} < \sigma^n \leq P \quad \text{or} \quad \sum_{v=i}^n q_v \cdot \sigma^{v-i} < Q \quad \text{for } i \geq 1.$$

Moreover, $\gcd(P, kQ) = 1$ such that (2.11) and (2.13) imply

$$(2.14) \quad \alpha_i \neq 0 \quad \text{and} \quad \beta_i \neq 0 \quad \text{for } i = 1 \cdots n.$$

Let $\iota_i \in \mathbb{R}^{n+1, n+1}$ be a matrix with 1 in the i th upper diagonal and 0 elsewhere such that

$$(2.15) \quad \iota_i \cdot s = (\sigma^{n-i}, \dots, \sigma, 1, 0, \dots, 0)^t \in \mathbb{R}^{n+1},$$

using s from (2.3). Then we are ready to describe C^{-1} as follows.

LEMMA 3. *The inverse of $C = C_n$ defined by (1.6) is given by*

$$(2.16) \quad \left[\begin{array}{c|c|c|c} P \cdot S & -k \cdot Q \cdot S & B & k \cdot A \\ \hline -Q \cdot S & P \cdot S & A & B \end{array} \right] \begin{array}{l} 0 \\ n+1 \\ n+2 \\ 2n+1 \end{array}$$

$0 \qquad 1 \qquad 2 \qquad n+1 \quad n+2 \quad 2n+1$

with

$$A := (\alpha_n s, \dots, \alpha_1 s) \in \mathbb{R}^{n+1, n},$$

and

$$B := ((\beta_n I + \iota_n) \cdot s, \dots, (\beta_1 I + \iota_1) \cdot s) \in \mathbb{R}^{n+1, n}.$$

Proof. For the matrices $A = (a_{ij})$ and $B = (b_{ij})$, we have

$$(2.17) \quad \begin{aligned} a_{ij} &= \alpha_{n-j+1} \cdot \sigma^{n-i}, \quad \text{and} \\ b_{ij} &= \begin{cases} \beta_{n-j+1} \cdot \sigma^{n-i}, & j \leq i, \\ \beta_{n-j+1} \cdot \sigma^{n-i} + \sigma^{j-i+1} & j \geq i+1 \end{cases} \end{aligned}$$

for $i = 0 \dots n, j = 1 \dots n$ (the row indices start with 0, the column indices with 1). Denote the matrix defined by (2.16) by Γ . Then for $0 \leq i, j \leq n$, we have

$$(\Gamma \cdot C)_{ij} = P \cdot s_i \cdot p_{n-j} - k \cdot Q \cdot s_i \cdot q_{n-j} + b_{i,j+1} - \sigma \cdot b_{ij}$$

where the third summand cancels for $j = n$, the fourth for $j = 0$. Using $\beta_0 = \beta_{n+1} = 0$ and (2.17) yields

$$(\Gamma \cdot C)_{ij} = \begin{cases} t(i, j) & \text{for } j < i, \\ t(i, j) + \sigma^{j-i} & \text{for } j = i, \\ t(i, j) + \sigma^{j-i} + \sigma^{j-i-1} & \text{for } j > i \end{cases}$$

using the abbreviation

$$t(i, j) := \sigma^{n-i} \cdot (P \cdot p_{n-j} - k \cdot Q \cdot q_{n-j} + \beta_{n-j} - \sigma \cdot \beta_{n-j+1}).$$

Therefore, for $0 \leq i, j \leq n$,

$$(2.18) \quad (\Gamma \cdot C)_{ij} = \sigma^{n-i} \cdot (P \cdot p_{n-j} - k \cdot Q \cdot q_{n-j} + \beta_{n-j} - \sigma \cdot \beta_{n-j+1}) + \delta_{ij}$$

using Kronecker's delta. Since later on we will need $|C^{-1}| \cdot |C|$, we write down the explicit formulae for the other components of $\Gamma \cdot C$. For $0 \leq i \leq n, n+1 \leq j \leq 2n+1$ derives

$$(2.19) \quad (\Gamma \cdot C)_{ij} = \sigma^{n-i} \cdot k \cdot (P \cdot q_{n-j} - Q \cdot p_{n-j} + \alpha_{n-j} - \sigma \cdot \alpha_{n-j+1});$$

for $n+1 \leq i \leq 2n+1, 0 \leq j \leq n$ derives

$$(2.20) \quad (\Gamma \cdot C)_{ij} = \sigma^{n-i} \cdot (-Q \cdot p_{n-j} + P \cdot q_{n-j} + \alpha_{n-j} - \sigma \cdot \alpha_{n-j+1});$$

and for $n+1 \leq i, j \leq 2n+1$ derives

$$(2.21) \quad (\Gamma \cdot C)_{ij} = \sigma^{n-i} \cdot (-k \cdot Q \cdot q_{n-j} + P \cdot p_{n-j} + \beta_{n-j} - \sigma \cdot \beta_{n-j+1}) + \delta_{ij}.$$

The identities (2.10) and (2.12) prove $(\Gamma \cdot C)_{ij} = \delta_{ij}$. □

For the condition of C using the ∞ -norm and $\alpha_i \neq 0$,

$$(2.22) \quad \begin{aligned} \|C_n\|_\infty \cdot \|C_n^{-1}\|_\infty &> \left\{ \sum_{\nu=0}^n (p_\nu + k \cdot q_\nu) \right\} \cdot \{ \sigma^n \cdot (P + k \cdot Q) \} \\ &= \left\{ \sum_{\nu=0}^n (\sigma^\nu p_\nu + k \sigma^\nu q_\nu) \right\} \cdot (P + kQ) \geq (P + k \cdot Q)^2. \end{aligned}$$

We calculate the sensitivity $s_{ij}(C)$ according to (2.1) for $0 \leq i \leq n, j = 0$. By (2.18) we have

$$(|C^{-1}| \cdot |C|)_{i\nu} \geq \sigma^{n-i} \cdot (P \cdot |p_{n-\nu}| + k \cdot Q \cdot |q_{n-\nu}| + |\beta_{n-\nu}| + \sigma \cdot |\beta_{n-\nu+1}|),$$

for $0 \leq \nu \leq n$ and by (2.19) we have

$$(|C^{-1}| \cdot |C|)_{i\nu} \geq \sigma^{n-i} \cdot k \cdot (P \cdot |q_{n-\nu}| + Q \cdot |p_{n-\nu}| + |\alpha_{n-\nu}| + \sigma \cdot |\alpha_{n-\nu+1}|),$$

for $n + 1 \leq \nu \leq 2n + 1$.

Using $\alpha_\nu, \beta_\nu \neq 0$ we get, for $0 \leq i \leq n$,

$$\begin{aligned} &(|C^{-1}| \cdot |C| \cdot |C^{-1}|)_{i0} \\ &= \sum_{\nu=0}^n (|C^{-1}| \cdot |C|)_{i\nu} \cdot |C^{-1}|_{\nu 0} + \sum_{\nu=n+1}^{2n+1} (|C^{-1}| \cdot |C|)_{i\nu} \cdot |C^{-1}|_{\nu 0} \\ &\geq \sigma^{n-i} \cdot \sum_{\nu=0}^n \{ (P \cdot |p_{n-\nu}| + k \cdot Q \cdot |q_{n-\nu}|) \cdot P \cdot \sigma^{n-\nu} + k \cdot (P \cdot |q_{n-\nu}| + Q \cdot |p_{n-\nu}|) \cdot Q \cdot \sigma^{n-\nu} \} \\ &\quad + \sigma^{n-i} \cdot \left\{ \sum_{\nu=0}^n (|\beta_{n-\nu}| + \sigma \cdot |\beta_{n-\nu+1}|) \cdot P \cdot \sigma^{n-\nu} \right. \\ &\quad \left. + \sum_{\nu=0}^n (|\alpha_{n-\nu}| + \sigma \cdot |\alpha_{n-\nu+1}|) \cdot kQ \sigma^{n-\nu} \right\} \\ &\geq \sigma^{n-i} \cdot P \cdot (P^2 + kQ^2 + kQ^2 + kQ^2) + \sigma^{n-i} \cdot P \cdot 4 \\ &= \sigma^{n-i} \cdot P \cdot (4P^2 - 3 + 4) > \sigma^{n-i} \cdot P \cdot (4P^2) \end{aligned}$$

using $k \cdot Q \geq P$. Together with $|C^{-1}|_{i0} = \sigma^{n-i} \cdot P \neq 0$,

$$S_{i0}(C) > 4P^2 \quad \text{for } 0 \leq i \leq n$$

follows. This proves the following theorem.

THEOREM 4. *The matrix C defined by (1.6) satisfies*

$$\|C\|_\infty \cdot \|C^{-1}\|_\infty > (P + k \cdot Q)^2$$

and there are components of C of which the sensitivity defined by (2.1) is greater than $4 \cdot P^2$.

3. Some examples. For given k , suitable pairs (P, Q) satisfying Pell's equation $P^2 - k \cdot Q^2 = 1$ are easily generated. Given some (P_0, Q_0) unequal, the trivial solution is $(1, 0)$, and successive solutions are

$$(P_{i+1}, Q_{i+1}) = (P_i P_0 + k Q_i Q_0, Q_i P_0 + P_i Q_0).$$

For a floating-point number system given by (1.1)–(1.3), a choice for σ is β^λ . Any expansion (1.5) of P, Q is suitable. The coefficients p_i, q_i are calculated successively.

Some bits can be saved by the following observation. If some coefficient p_i is divisible by β or by a power of β , then p_i and the following $p_j, j > i$ are expressed with a corresponding exponent. If the last digit m_λ in the mantissa of p_{i+1} is equal to $\beta - 1$, then p_i can be replaced by $p_i - \sigma$ and p_{i+1} by $p_{i+1} + 1$, the latter being divisible by β .

For example, let $P = 73942, \beta = 10, \sigma = 100$. Then expanding P yields $(p_2, p_1, p_0) = (7, 39, 42)$ and this is reduced by the method just described to $(p_1, p_0) = (74 \cdot 10^1, -58)$. This method is especially useful for base 2.

For a given number P , the corresponding coefficients $p_i, i = 0 \dots n$ can be calculated by the following algorithm:

```

e = 0; i = 0;
repeat
  while P mod beta = 0 do { P = P/beta; e = e+1 };
  q = floor(P/beta); r = P - q*beta;
  if (q mod beta != beta-1) or (q < beta)
    then { p_j = r * beta^e; P = q }
    else { p_i = (r - sigma) * beta^e; P = q + 1 };
  i = i + 1
until P = 0;
    
```

For $k = 2$, successive pairs P, Q are $(3, 2), (17, 12), (99, 70) \dots$. In Table 1 we display some values for p_i, q_i for single and double precision. For the individual value of n (resulting in a $2n \times 2n$ -matrix C) we choose the maximum values (P, Q) being representable by (p_{n-1}, \dots, p_0) and (q_{n-1}, \dots, q_0) . In the columns of Table 1, the condition number is given followed by the coefficients p_i and q_i , both in descending order. The coefficients are given by two numbers m and e such that $m \cdot 2^e$ is the actual coefficient. For example, $q_4 = 1175 \cdot 2^{22}$ for $n = 5$ (yielding a 10×10 -matrix). Our algorithm yields a higher condition than the expected maximum $4 \cdot 2^{24 \cdot 2n} \approx 7 \cdot 10^{72}$, especially for this 10×10 -matrix.

For double precision we choose different values for k yielding the coefficients in Table 2. These coefficients are, of course, only samples used to construct matrices of the general form (1.6). We conclude by writing the 6×6 -matrix for single precision explicitly.

TABLE 1
 p_i, q_i for binary format, 24 bit precision; $k = 2$.

Cond	1.3E+030	2.2E+044	6.5E+060	1.1E+078	4.8E+090	1.7E+107
p_i	15248163 2	3527199 3	6929233 6	425393 14	2161033 8	8490761 10
	11171905 0	6746489 1	9763077 3	6127903 11	5075327 7	15520103 6
q_i		-8816797 0	12608263 1	-10707825 7	8241033 6	6855055 5
	84235 9		-6160127 0	7194379 1	-9934673 5	6997339 4
	-3559681 3	1247053 4		-2285085 0	-5752371 1	-11831695 3
		13508351 2	1224927 8		12291875 0	9051609 1
		-14061827 1	-5131195 6	1175 22		-11093871 0
			14870387 5	-14199789 15	47753 13	
			-7145793 4	12492253 13	-15523515 12	3001937 11
				9093109 10	-1620555 9	12103369 10
				10074835 1	14867027 6	-13213329 9
					14366575 3	-9497253 7
				-4879973 1	-3241495 4	
					8507481 3	
					-1367575 2	

TABLE 2
 p_i, q_i for binary format, 53 bit precision.

Cond k	7.0E+066 32	3.4E+097 2	2.1E+131 32	1.4E+164 2
p_i	8384758637032543 5 -3529290569461695 0	119071610094027 9 -3183251058136493 3 -8183182949466111 0	1838140087490775 8 -6618243915631817 2 -7698164339527309 1	1217131843483323 9 5555590710757647 8 -1048381871128883 4
q_i	5928919690858185 3 -6097772977423311 1	84196342944287 9 891386017353869 8 -1900818942150157 7	162470165079445 9 6774769086897599 6 4831599480133437 3 -5900891544265983 0	4113071334050663 3 -3228782923936605 0 1721284360250283 8 292142371452983 6 -4351444206118847 4 1403045714199203 2 -2787903664869301 1

It is exactly storable with only 24 bits in the mantissa (and therefore in almost any floating-point number system) but matrix inversion will "fail" in almost any floating-point format available because, due to the condition number $2.2 \cdot 10^{44}$, an equivalent of approximately 44 decimal digits precision would be necessary:

$$\begin{pmatrix} 3527199 \cdot 2^3 & 6746489 \cdot 2^1 & -8816797 \cdot 2^0 & 1247053 \cdot 2^5 & 13508351 \cdot 2^3 & -14061827 \cdot 2^2 \\ 1247053 \cdot 2^4 & 13508351 \cdot 2^2 & -14061827 \cdot 2^1 & 3527199 \cdot 2^3 & 6746489 \cdot 2^1 & -8816797 \cdot 2^0 \\ 1 & -2^{24} & & & & \\ & 1 & -2^{24} & & & \\ & & & 1 & -2^{24} & \\ & & & & 1 & -2^{24} \end{pmatrix}.$$

To generate this matrix, the values $P = 7942546277405390632803$ and $Q = 5616228332641321147898$ have been used.

MATLAB [2] delivers as an estimation for the condition number of the matrix the (almost) correct answer ∞ .

REFERENCES

- [1] G. H. HARDY AND E. WRIGHT, *An Introduction to the Theory of Numbers*, Fifth Edition, Oxford Science Publications, Oxford, 1980, 1981, p. 442.
- [2] PRO-MATLAB *User's Guide*, Vers. 32-SUN, The MathWorks, Inc., Sherborn, MA, 1987.
- [3] J. ROHN, *New condition numbers for matrices and linear systems*, Computing, 41 (1989), pp. 167-169.