

Effects of Underflow on Solving Linear Systems

James Demmel

Computer Science Division
University of California
Berkeley, CA

ABSTRACT

Software to solve systems of linear equations by Gaussian elimination has in the past ignored the effects of underflow. But when underflows are replaced by zeroes, this software can give spurious though plausible results much worse than could be blamed on roundoff. When underflow is gradual, as in the proposed IEEE standard for floating point arithmetic, the same software gives provably more reliable results. To achieve the same reliability without gradual underflow, but with underflows set to zero, complicated tests must be inserted into the software.

1. Introduction

In this paper we examine the effects of underflow on solving systems of linear equations using Gaussian elimination. In particular, we contrast the effects of gradual underflow¹ and "store zero" on the accuracy and stability of the algorithm.

Our ultimate goal is to decide whether reliable software for solving linear systems in the presence of underflow can be written at a reasonable cost, where by reliable we mean a piece of software that ideally

1. produces accurate results whenever they can be represented,
2. gives a warning whenever the computed results are inaccurate, and
3. avoids giving warnings when it is possible to compute accurate results.

Our motivation for this study is twofold. First, that underflow is a problem at all in linear equation solvers is not generally recognized. For example, the LINPACK Users' Guide² states in its introduction that "Any underflows which occur are harmless." This statement is almost always true, since underflow's status as a borderline phenomenon means it may be ignored safely in most situations, but not always. The advantage of gradual underflow over store zero is that a cleaner line can be drawn between problems where underflow can be ignored

and those where it cannot, thus helping the writer of reliable software who wants to be able to guarantee when his program will work. Our first motivation, then, is to make people aware that underflow can give reasonable looking but totally inaccurate results if it is handled poorly or ignored.

Second, we wish to attenuate the controversy surrounding the recent decision by the IEEE Microprocessor Standards Committee on how to handle underflow. The decision to use gradual underflow instead of the usual store zero approach came after much argument about which one made reliable numerical software easier to write¹. In this paper we show that gradual underflow makes writing reliable linear system solvers significantly easier than store zero.

Specifically, the algorithm to solve $Ax=b$ is as follows:

- (1) Decompose $A=LU$ (L lower triangular, U upper triangular) using pivoting, so that the diagonal of L contains all 1's and no entries of L exceed 1 in absolute value.
- (2) Solve $Ly=b$ for y (forward substitution).
- (3) Solve $Ux=y$ for x (back substitution).

The only gradual underflows that can possibly contribute significantly to the residual turn out to be underflows in the final solution vector x . Intermediate gradual underflows will be shown to contribute an error with a bound scarcely worse than the bound for the error contributed by roundoff alone. Thus, by using gradual underflow we are able to satisfy the second two goals of reliable software since an alarm need not be raised unless the results themselves underflow and are not representable (in which case the data would have to be scaled to avoid underflow).

In contrast, storing zero in place of intermediate underflows during any stage of solution can introduce significant errors, possibly producing reasonable looking results whose error greatly exceeds the uncertainty attributable to roundoff alone (see the Examples).

The second and third goals of reliable software can be achieved, but not easily, with store zero. We must either do tedious testing to tell which underflows require a warning, or naively raise an alarm whenever an underflow occurs. This naive approach produces "paranoid" code and many false alarms. To tell which small class of possible gradual underflows can contribute significant errors and require a warning is easy: they are the underflows in the answer itself.

It is important to understand how the unavoidable uncertainty due to roundoff can affect our ability to meet our first goal: produce accurate results. Guaranteeing

results x to within a certain accuracy is not possible for a price everybody is willing to pay. Unless we pay the price, by doing iterative refinement for example (discussion of iterative refinement is beyond the scope of this paper), we can only guarantee that our computed answer x is the solution of a new unknown problem perturbed slightly from our original problem. This is the nature of Gaussian elimination, and confirmed by backwards error analysis, the approach used in this paper; we want to bound that perturbation, to bound how different the new problem is from the old. If the problem is ill-conditioned, then the solution x of the new problem may be very different from the solution x of the original problem, and there is no way, at negligible cost in time and storage, to tell. Thus, our original goal of producing accurate results should be modified to read: "satisfy the equations as closely as possible; achieve a tiny residual $Ax - b$." This is an achievable goal. (Discussion of iterative refinement is beyond the scope of this paper.)

In the remainder of this paper we will present examples to demonstrate the typical effects of underflow, describe the model of arithmetic used and the approach used in the error analysis, and present the conclusions drawn from the error analysis.

2. Examples

We use G.U. to denote arithmetic performed with gradual underflow, and S.Z. to denote that performed with underflows set to zero; we assume the reader is acquainted with both kinds of arithmetic^{1,3,4}. In particular, we perform G.U. in the "normalizing mode" rather than the "warning mode". We will present our model of arithmetic briefly here, and give a more precise model of rounding and underflow errors in the next section.

Let ϵ be the rounding error of the arithmetic, and λ be the underflow threshold. Suppose, for example, a binary floating point number is represented as $f \cdot 2^e$ where f lies between 1 and 2 with an n bit fraction, and e is an integer exponent in the range $e_{\min} \leq e \leq e_{\max}$. Then $\epsilon = 2^{-n}$ (the difference between 1 and the next larger number). The underflow threshold is $\lambda = 2^{e_{\min}}$; this is the smallest positive number in S.Z. arithmetic, and the smallest normalized number for G.U. Whenever a nonzero number smaller than λ is generated by arithmetic, underflow is signalled and something special happens; S.Z. replaces the number by zero, and G.U. by a nearest "denormalized" number or zero. Denormalized numbers in G.U. arithmetic consist of an arithmetic progression from 0 to λ with common separation equal to $\mu = \epsilon\lambda$.

We denote the usual condition number of the matrix A by

$$k(A) = \|A\|_{\infty} \cdot \|A^{-1}\|_{\infty}$$

and a new set of condition numbers by

$$\text{Cond}(A, x) = \frac{\| |A^{-1}| |A| |x| \|_{\infty}}{\|x\|_{\infty}}$$

$$\text{Cond}(A) = \| |A^{-1}| |A| \|_{\infty}$$

These new condition numbers, due to Skeel⁵, will be discussed more fully in section 3 below. Note $\text{Cond}(A) \geq \text{Cond}(A, x)$ for all x .

In this section we present four examples of the effects of underflow on performing Gaussian elimination. The first example shows how store zero can produce a reasonable looking but completely inaccurate decomposition of a well conditioned matrix, whereas gradual underflow

either produces the correct decomposition or correctly decides the matrix is singular. (There are no rounding errors nor pivot growth in this example.) The second example shows that G.U. produces the correct decomposition of a well conditioned matrix which S.Z. incorrectly decides is singular. Third, we present an innocuous looking ordinary differential equation and show that the linear system arising from trying to solve it numerically leads to underflow which is handled correctly by G.U. and not by S.Z. Finally, we present an example which shows that regardless of whether we use G.U. or S.Z., Gaussian elimination can only guarantee small residuals, not an accurate answer, even when the matrix A is well-conditioned in the sense that $\text{Cond}(A)$ is small.

2.1 Example 1

Consider the family of matrices $A(x)$ where

$$A(x) = \lambda \cdot \begin{bmatrix} 2 & & 1 \\ & 2 & 1 \\ & & 2 & 1 \\ & & & 2 & 1 \\ 1 & 1 & 1 & 1 & x \end{bmatrix} \quad (1)$$

(blanks denote zero entries). The LU decomposition obtained by G.U. is

$$L^{G.U.}(x) \cdot U^{G.U.}(x) = \begin{bmatrix} 1 & & & & \\ & 1 & & & \\ & & 1 & & \\ & & & 1 & \\ .5 & .5 & .5 & .5 & 1 \end{bmatrix} \cdot \lambda \cdot \begin{bmatrix} 2 & & 1 \\ & 2 & 1 \\ & & 2 & 1 \\ & & & 2 & 1 \\ & & & & x-2 \end{bmatrix} \quad (2)$$

$= A(x)$ exactly, and by S.Z. is

$$L^{S.Z.}(x) \cdot U^{S.Z.}(x) = \begin{bmatrix} 1 & & & & \\ & 1 & & & \\ & & 1 & & \\ & & & 1 & \\ .5 & .5 & .5 & .5 & 1 \end{bmatrix} \cdot \lambda \cdot \begin{bmatrix} 2 & & 1 \\ & 2 & 1 \\ & & 2 & 1 \\ & & & 2 & 1 \\ & & & & x \end{bmatrix} \quad (3)$$

$= A(x) + E$, where the error matrix E equals

$$E = \lambda \cdot \begin{bmatrix} & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & -2 \end{bmatrix} \quad (4)$$

We see S.Z. causes a relatively large error in the $U(x)_{55}$ entry, whereas G.U. gives the correct decomposition. When $x=2$, using S.Z. leads us to conclude that the matrix is far from singular, when in fact it is exactly singular. Note that the matrix $A(x)$ is well conditioned when x is far from 2, and if x is a small integer no rounding errors occur in either decomposition.

2.2 Example 2

Let

$$A = \begin{bmatrix} 2\lambda & 3\lambda \\ \lambda & 2\lambda \end{bmatrix} \quad (5)$$

a well conditioned matrix. Using G.U. we obtain

$$L^{G.U.} \cdot U^{G.U.} = \begin{bmatrix} 1 & \\ .5 & 1 \end{bmatrix} \cdot \begin{bmatrix} 2\lambda & 3\lambda \\ \lambda & \lambda/2 \end{bmatrix} = A \quad (6)$$

but by using S.Z. we obtain

$$L^{S.Z.} \cdot U^{S.Z.} = \begin{bmatrix} 1 & & \\ 5 & 1 & \\ & & \ddots \end{bmatrix} \cdot \begin{bmatrix} 2\lambda & 3\lambda \\ & 0 \\ & & \ddots \end{bmatrix}. \quad (7)$$

Thus, G.U. correctly decomposes the matrix A , whereas S.Z. incorrectly makes the matrix look singular.

2.3 Example 3

Consider the ordinary differential equation

$$\dot{x}(t) = \frac{1-(t/T)^M}{T-t} x(t), \quad x(T_0) = c. \quad (8)$$

We try to solve this equation numerically by replacing $x(t)$ by the truncated power series $\sum_{n=1}^N x_n t^n$, the function $(1-(t/T)^M)/(T-t)$ by its (finite) power series, and then equating coefficients of equal powers of t on both sides of equation (8). After we scale the last row (which represents the initial condition) down to have largest entry equal to 1, we get the linear system $Ax=b$, where

$$A = \begin{bmatrix} N & -1/T & -1/T^2 & \dots & -1/T^N \\ N-1 & -1/T & \dots & \dots & -1/T^{N-1} \\ & N-2 & \dots & \dots & -1/T^{N-2} \\ & & & & \vdots \\ & & & & \vdots \\ & & & & 1 & -1/T \\ 1 & 1/T_0 & 1/T_0^2 & \dots & 1/T_0^{N-1} & 1/T_0^N \end{bmatrix}. \quad (9)$$

$b^T = (0, \dots, 0, c/T_0^N)$, and $x^T = (x_N, \dots, x_0)$.

We chose $M=15$, $N=14$, $T=512$, $T_0=500$, and $c=100$ for this example. We used an implementation of the IEEE Floating Point Standard⁴ on a VAX 11/780⁶. ϵ was 2^{-23} and μ was 2^{-149} . There was a switch on the compiler to enable/disable G.U., so we were able to obtain numerical results using both G.U. and S.Z.

L and U have a simple structure. L will be zero below the diagonal, except for the last row, which is graded from $L_{15,1} \approx 7.14285_{10}^{-2}$ down to $L_{15,14} \approx 5.34728_{10}^{-35}$. U is identical to A in all but its last row.

$$L = \begin{bmatrix} 1 & & & & \\ 0 & 1 & & & \\ 0 & 0 & & & \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & \dots & 1 & \\ L_{15,1} & L_{15,2} & \dots & L_{15,14} & 1 \end{bmatrix} \quad (10)$$

$$U = \begin{bmatrix} N & -1/T & -1/T^2 & \dots & -1/T^N \\ N-1 & -1/T & \dots & \dots & -1/T^{N-1} \\ & N-2 & \dots & \dots & -1/T^{N-2} \\ & & & & \vdots \\ & & & & \vdots \\ & & & & 1 & -1/T \\ & & & & & U_{15,15} \end{bmatrix} \quad (11)$$

A 's columns are badly scaled, although this is not obvious because no row nor column is drastically smaller in norm than any other; nonetheless, bad scaling causes A to appear very ill conditioned, and this ill conditioning shows up in the last row of U , making $U_{15,15}$ very small,

barely above the underflow threshold. S.Z. and G.U. compute all elements of L and U identically except for $U_{15,15}$. In fact, all additions in the computation of L add normalized numbers with like magnitudes and like signs, so no cancellation, loss of significance, nor underflows occur. If the exponent range were unbounded, so underflow never happened, the correct value $U_{15,15} \approx 2.09281_{10}^{-37}$ would be computed. This is the value computed using G.U. But when S.Z. is used instead, the computed value is $U_{15,15}^{S.Z.} \approx 1.72783_{10}^{-37}$, a relative difference of .174 from the correct value. All additions in the computation of $U_{15,15}$ involve numbers of like magnitude and sign, so cancellation cannot be blamed for the discrepancy. This relative difference in the last entry of U is very important, because one divides by $U_{15,15}$ in the course of solution. Thus, the computed solution $x^{G.U.}$ is very close to the true x , and the relative difference in solution vectors is

$$\frac{\|x^{G.U.} - x^{S.Z.}\|_{\infty}}{\|x^{G.U.}\|_{\infty}} \approx .211.$$

Thus, G.U. obtains markedly better results than S.Z. This example is very interesting because there is nothing obviously wrong with the matrix. All its entries are unexceptional normalized numbers, and every row and every column contains at least one number no tinier than $1/T \approx .00195$ and none larger than $N=14$, yet 11 out of 14 products $L_{15,j} \cdot U_{j,15}$ in the sum for $U_{15,15}$ underflow just slightly below the underflow threshold. Since the true value of $U_{15,15}$ is itself not much larger than the underflow threshold, this makes for a large relative error.

This example was chosen to be simple and realistic; even though it can be solved analytically, it could be changed easily into a two dimensional problem without an explicit solution, but with the same sensitivity to underflow.

We repeat that even though A appears very ill conditioned, since $k(A) \approx 1/\lambda$ (i.e. near the overflow threshold in most arithmetics), it is also well conditioned in the sense that $\text{Cond}(A, x) \approx 5.5$. We will discuss the significance of this example later in section 4.

2.4 Example 4

Let

$$A = \begin{bmatrix} G & G \\ g & 2g \end{bmatrix}, \quad A^{-1} = \begin{bmatrix} 2/G & -1/g \\ -1/G & 1/g \end{bmatrix},$$

where g/G underflows to 0 using either S.Z. or G.U. The L obtained is thus the identity matrix since $L_{2,1} = fl(g/G) = 0$, and so the L and U obtained are the exact factors of the matrix

$$A+E = \begin{bmatrix} G & G \\ 0 & 2g \end{bmatrix},$$

which is a very different matrix than A . If $b^T = (G, 0)$, then $x = A^{-1}b = (2, -1)^T$, whereas $\hat{x} = (A+E)^{-1}b = (1, 0)^T$, so \hat{x} does not resemble x at all. The residual r is however guaranteed to be small, in the sense that $\|r\|_{\infty} / (\|A\|_{\infty} \|\hat{x}\|_{\infty} + \|b\|_{\infty})$ is small:

$$\frac{\|r\|_{\infty}}{\|A\|_{\infty} \|\hat{x}\|_{\infty} + \|b\|_{\infty}} = \frac{\|E\hat{x}\|_{\infty}}{\|A\|_{\infty} \|\hat{x}\|_{\infty} + \|b\|_{\infty}} \leq \frac{g|\hat{x}_1|}{G|\hat{x}_1| + G|\hat{x}_2|} \leq \frac{g}{G} \leq \lambda\epsilon/2.$$

Of course A is an exceedingly ill conditioned matrix in the sense that $k(A) \approx 2G/g$ is beyond the reciprocal of the underflow threshold, so we would be inclined not to trust our results anyway. However, $\text{Cond}(A)$ is only 7. This is true because $\text{Cond}(A) = \text{Cond}(DA)$ for any non-singular diagonal matrix D , so A has the same condition number as the utterly tame matrix

$$\begin{bmatrix} G^{-1} & \\ & g^{-1} \end{bmatrix} A = \begin{bmatrix} 1 & 1 \\ 1 & 2 \end{bmatrix}$$

Needless to say, in the absence of underflow we would compute a very accurate solution. We will return to this example later to explain why we can get inaccurate results from a matrix with a small condition number $\text{Cond}(A)$.

3. Error Analysis

3.1 Assumptions

In this section we define our notation and our model of arithmetic. ϵ will denote the level of roundoff error and μ the smallest nonzero number when using G.U. Then $\lambda = \mu/\epsilon$ will be the underflow threshold and smallest normalized number. Using S.Z. the only number smaller than λ in magnitude is zero; using G.U. the numbers between λ and zero in magnitude are called denormalized numbers. Zero and all numbers no smaller than λ in magnitude are normalized in both G.U. and S.Z.

Let \bullet be one of the operations $(+, -, *, /)$ and let $\mathfrak{fl}(a \bullet b)$ denote the floating point result of the indicated computation. Traditionally, error analyses have used the formula

$$\mathfrak{fl}(a \bullet b) = (a \bullet b)(1 + \epsilon) \quad (12)$$

unless $a \bullet b$ underflows or overflows. To take underflow into account, we write⁹

$$\mathfrak{fl}(a \bullet b) = (a \bullet b)(1 + \epsilon) + \eta \quad (13)$$

unless $a \bullet b$ overflows. In the case of G.U. we have the following constraints on ϵ and η :

- (1) $|\epsilon| \leq \epsilon$ and $|\eta| \leq \lambda \epsilon$,
- (2) $\eta \epsilon = 0$,
- (3) $\eta = 0$ if \bullet is either addition or subtraction.

In the case of S.Z. we have the following somewhat different constraints on ϵ and η :

- (1) $|\epsilon| \leq \epsilon$ and $|\eta| \leq \lambda$, and
- (2) $\eta \epsilon = 0$.

Note that η need not be 0 when performing an addition or subtraction using S.Z.

We assume no overflow occurs.

We ignore terms which are $O(\epsilon^2)$ or $O(\epsilon\mu)$. This allows us to replace $1/(1+\epsilon_1)$ by $1-\epsilon_2$ (where $|\epsilon_1| \leq \epsilon$), or $\eta_1/(1+\epsilon)$ by η_2 (where $|\eta_1| \leq \lambda$ for S.Z. or $\lambda\epsilon$ for G.U.) when convenient.

By replacing every appearance of an original datum a by $1 \cdot a$, and using the above formula for error in multiplication, we can take into account effects of rounding and underflow errors in the original data.

$\|A\|$ and $\|b\|$ denote matrix and vector norms. $\|b\|_\infty$ denotes the infinity norm of b , namely $\max_j |b_j|$, and

$\|A\|_\infty$ denotes the infinity norm of A , $\|A\|_\infty = \max_{i \neq 0} \|Ax\|_\infty / \|x\|_\infty = \max_j \sum_i |A_{ij}|$.

Then $k(A) = \|A\|_\infty \|A^{-1}\|_\infty$ denotes usual the condition number of the matrix A .

$|A|$ ($|b|$) denotes the matrix (vector) whose entries are the absolute values of the corresponding entries in A (b). Inequalities like $|A| > |B|$ ($|a| > |b|$) are meant componentwise.

3.2 Approach

As stated in the introduction, we use backward error analysis. Thus, when Gaussian elimination is used to solve

$$Ax = b \quad (14)$$

for x it generates instead an approximation $\hat{x} = x + \delta x$ which satisfies some perturbed problem

$$(A + \delta A)\hat{x} = b + \delta b \quad (15)$$

The task of backwards error analysis is to infer bounds on δA and δb from the details of the arithmetic used to implement the elimination process. These bounds can be used in turn to bound the residual

$$r = A\hat{x} - b = -\delta A\hat{x} + \delta b = A\delta x \quad (16)$$

and then the error δx .

We proceed by breaking the elimination process into three steps:

- 1) Decomposition : Try to decompose A into LU and obtain L and U where $A=LU - \Delta A$; bound ΔA in terms of A .
- 2) Forward Substitution : Try to solve $Ly=b$ and obtain \hat{y} where $(L+\delta L)\hat{y} = b + \Delta b$; bound δL in terms of L and Δb in terms of b .
- 3) Backward Substitution : Try to solve $Ux=\hat{y}$ and obtain \hat{x} where $(U+\delta U)\hat{x} = \hat{y} + \delta \hat{y}$; bound δU in terms of U and $\delta \hat{y}$ in terms of \hat{y} .

Thus, combining the above three steps we obtain

$$(A + \Delta A + L\delta U + \delta LU + \delta L\delta U)\hat{x} = b + \Delta b + (L + \delta L)\delta \hat{y} \quad (17)$$

The residual $r = A\hat{x} - b$ is

$$r = -\delta A\hat{x} + \delta b = -(\Delta A + L\delta U + \delta L(U + \delta U))\hat{x} + \Delta b + (L + \delta L)\delta \hat{y} \quad (18)$$

Equation (18) shows that we may choose δA and δb to be

$$\delta A = \Delta A + L\delta U + \delta L(U + \delta U) \quad (19)$$

and

$$\delta b = \Delta b + (L + \delta L)\delta \hat{y} \quad (20)$$

The usual error analyses attribute the various perturbations ΔA , δL , Δb , δU , and $\delta \hat{y}$ directly to roundoff errors during the elimination process. Underflow contributes additional perturbations thus:

- (1) during decomposition \dots to ΔA and δL ,
- (2) during forward substitution \dots to Δb ,
- (3) during backward substitution \dots to $\delta \hat{y}$ and \hat{x} .

These additional perturbations, of the order of λ for S.Z. and $\mu = \lambda\epsilon$ for G.U., all turn out to be comparable with or less than the respective terms to which they contribute. Our task is to summarize quantitatively the perturbations' cumulative effect.

Wilkinson's approach⁷ is to determine a bound ω_w on the errors

$$\|\delta A\|_\infty \leq \omega_w \|A\|_\infty \quad \text{and} \quad \|\delta b\|_\infty \leq \omega_w \|b\|_\infty \quad (21)$$

whence

$$\|r\|_{\infty} \leq \omega_w [\|A\|_{\infty} \|\hat{x}\|_{\infty} + \|b\|_{\infty}] \quad (22)$$

and then it will follow that the error δx is bounded:

$$\frac{\|\delta x\|_{\infty}}{\|x\|_{\infty} + \|\hat{x}\|_{\infty}} \leq \omega_w k(A) \quad (23)$$

The detailed derivation of ω_w from the details of the arithmetic is given elsewhere⁹. Theorem 1 below states simple requirements on A and b that ensure ω_w will be scarcely worse if underflow occurs than if it does not.

Skeel's approach⁶, modified slightly here, is to determine a bound ω_s on the relative error in each entry of A and b :

$$|\delta A| \leq \omega_s |A| \quad \text{and} \quad |\delta b| \leq \omega_s |b| \quad (24)$$

From these inequalities follows a bound upon the error δx :

$$\frac{\|\delta x\|_{\infty}}{\|x\|_{\infty}} \leq \omega_s \cdot \frac{\| |A^{-1}| |A| |x| + |A^{-1}| |b| \|_{\infty}}{(1 - \omega_s \| |A^{-1}| |A| \|_{\infty}) \|x\|_{\infty}} \quad (25)$$

(provided the denominator is positive.) This motivates defining the following condition numbers:

$$\text{Cond}(A, x) = \frac{\| |A^{-1}| |A| |x| \|_{\infty}}{\|x\|_{\infty}} \quad (26a)$$

$$\text{Cond}(A) = \| |A^{-1}| |A| \|_{\infty} \quad (26b)$$

$\text{Cond}(A)$ is an upper bound for $\text{Cond}(A, x)$ for all x ; the error bounds are useful only if $\omega_s \text{Cond}(A) < 1$.

The bound ω_s in (24) is not inferred directly from equations (19) and (20). Instead, following Oetti and Prager⁸ and Skeel⁶ we use an expression for ω_s obtainable from (16) in terms of the residual r :

$$\omega_s = \max_i \frac{|\tau_i|}{(|A| |\hat{x}| + |b|)_i} \quad (27)$$

where the max is over those i for which the denominator is nonzero. Following Skeel, we overestimate ω_s by analyzing the elimination process to infer an inequality

$$\|r\|_{\infty} \leq \omega'_s \| |A| |\hat{x}| + |b| \|_{\infty} \quad (28)$$

from which we compute the overestimate $\bar{\omega}_s$ as

$$\bar{\omega}_s = \frac{\max_i (|A| |\hat{x}| + |b|)_i}{\min_i (|A| |\hat{x}| + |b|)_i} \cdot \omega'_s \quad (29)$$

(where the min in the denominator is over the nonzero values of $(|A| |\hat{x}|)_i$ only). Unfortunately $\bar{\omega}_s$ can be a gross overestimate of ω_s , as we will see when we return to example 3 later.

The detailed derivation of ω'_s is given in reference 9. Theorem 2 below states requirements on A and b that ensure ω'_s will be scarcely worse if underflow occurs than if it does not. These requirements on A and b are nearly identical to the requirements in the Wilkinson style analysis.

3.3 Results

Theorem 1: Wilkinson style error analysis of solving $Ax=b$ in the presence of underflow

Let $\alpha_{\max} = \max_j |A_j|$, and $g = [\text{largest intermediate result appearing in the decomposition}] / \alpha_{\max}$. g is the "pivot growth factor" and is $\leq 2^{n-1}$.

Our objective is to estimate a bound ω_w for which

$$\|r\|_{\infty} \leq \omega_w [\|A\|_{\infty} \|\hat{x}\|_{\infty} + \|b\|_{\infty}] \quad (22)$$

In the absence of underflow, we can use

$$\omega_w = n^3 \epsilon g / 2 \quad (30)$$

If underflow occurs then

$$\omega_w = 3 n^3 \epsilon g / 2 \quad (31)$$

provided certain conditions are met.

For G.U. these conditions are:

$$g \alpha_{\max} \geq \lambda \quad \text{if there are any underflows during triangular decomposition}$$

$$\|b\|_{\infty} \geq \frac{\lambda}{n} \quad \text{if there are any intermediate underflows during forward and back substitutions} \quad (32)$$

$$\frac{\|b\|_{\infty}}{\alpha_{\max}} \geq \frac{2\lambda}{n^2} \quad \text{if the solution } \hat{x} \text{ itself underflows.}$$

For S.Z. the above conditions still apply but λ must be increased to λ/ϵ .

Proof: See reference 9.

Theorem 2: Skeel style error analysis of solving $Ax=b$ in the presence of underflow

Let $\alpha_j = \max_i |A_{ij}|$, and $g_c = \max_j (\text{largest intermediate result appearing in the decomposition in column } j) / \alpha_j$.

g_c is the "columnwise pivot growth factor" and is $\leq 2^{n-1}$.

Our objective is to estimate a bound ω'_s for which

$$\|r\|_{\infty} \leq \omega'_s \| |A| |\hat{x}| + |b| \|_{\infty} \quad (28)$$

In the absence of underflow we can use

$$\omega'_s = n^3 \epsilon g_c \quad (33)$$

If underflow occurs, then

$$\omega'_s = 3 n^3 \epsilon g_c / 2 \quad (34)$$

provided certain conditions are met.

For G.U. these conditions are:

$$g_c \alpha_j \geq \lambda \quad \text{for all } j, \text{ if there are any underflows during triangular decomposition}$$

$$\|b\|_{\infty} \geq \frac{\lambda}{2n} \quad \text{if there are any intermediate underflows during forward and back substitutions} \quad (35)$$

$$\frac{\|b\|_{\infty}}{\alpha_{\max}} \geq \frac{\lambda}{n^2} \quad \text{if the solution } \hat{x} \text{ itself underflows.}$$

For S.Z. the above conditions apply with except λ must be increased to λ/ϵ .

Proof: See reference 9.

The theorems indicate how to write software that will solve $Ax=b$ reliably despite underflow, and how the requirements for G.U. differ from those for S.Z. To keep the residual small in the sense of a Wilkinson style error analysis, we appeal to Theorem 1. With G.U., as long as one normalized number appears during the decomposition ($g \alpha_{\max} \geq \lambda$), residual with underflow has a bound not much worse than residual without underflow. If there are intermediate underflows while solving the triangular systems, as long as some component of b is normalized ($\|b\|_{\infty} \geq \lambda$), residual with underflow has a bound

scarcely worse than without underflow. If the answer \hat{x} itself underflows, we can either issue an error message (which would be very reasonable since the first goal of reliable software is only to compute an answer if it is representable) or test to see if $\|b\|_{\infty}/\alpha_{\max}$ is not too small.

All these requirements are natural ones to make, since they say that when a problem's inputs and its computed solution are normalized numbers, we should expect the residual to be scarcely worse with underflow than without. Thus, the only gradual underflows which can cause concern in a problem with normalized inputs are underflows in the solution itself.

In contrast, the bounds for S.Z. are all higher by a factor of $1/\epsilon$. Thus, using S.Z. we can neither solve as many problems as with G.U., nor decide so easily which underflows matter. Thus, from the point of view of a Wilkinson style error analysis, G.U. makes writing reliable software easier.

Theorem 2 shows that Skeel style bounds for the residual are scarcely worse with underflow than without provided conditions are satisfied that are almost the same as in Theorem 1. Therefore the previous paragraphs' comments remain valid provided, when underflow is gradual, at least one normalized number appears in each column of A , rather than just somewhere in A , before or during the decomposition process.

The programs in Appendices 1 and 2 show the few extra lines of code needed to implement the test for gradual underflows in the solution \hat{x} in terms of the proposed IEEE standard.

4. Examples 3 and 4 Revisited

We wish to emphasize that we have only derived conditions under which *residual bounds* with underflow are about the same as without underflow. There is no way using this analysis to say how closely this bound will be approached with and without underflow, or how accurate the computed solution will be.

In example 4 above, the matrix A and vector b satisfy all the conditions of Theorems 1 and 2 for G.U. as well as S.Z., so the residual is small, but the answer \hat{x} is totally inaccurate. This inaccuracy can be explained either by the huge condition number $k(A)$ above overflow threshold, or the large backwards error in equation (27): $\omega_b = 1$. In this case ω_b 's upper bound $\bar{\omega}_b$ in equation (29) is also 1. Thus, having a small value of $\text{Cond}(A)$ is not sufficient to guarantee accuracy given a small residual ω_b ((28)), although a small value of $k(A)$ combined with a small residual ω_b is enough, as can be seen from (23).

Example 3 is another case where the conditions of Theorems 1 and 2 hold, but now G.U. successfully computes the last pivot $U_{15,15}$ and an accurate solution \hat{x} while S.Z. does not. Again, we have a problem where $k(A)$ is huge and $\text{Cond}(A, x)$ is small. Now the ω_b of equation (27) is $\approx 5.23_{10}^{-8}$, verifying the high accuracy of solution. Unfortunately the bad scaling of the matrix causes the upper bound $\bar{\omega}_b$ of equation (29) to be 2.0_{10}^{20} . This example demonstrates the occasionally intense pessimism of Skeel's approach.

In summary, the significance of examples 3 and 4 is to show that maintaining a small residual in the face of underflow does not guarantee an accurate solution \hat{x} , although we conjecture that for not terribly ill conditioned matrices G.U. will provide answers at least as accurate as provided by S.Z.

5. Conclusions

We have proven something quite unremarkable: if underflows are gradual, then we continue to get what we have come to expect from Gaussian elimination. That is, we get a small residual as long as the inputs and outputs are all representable (normalized) numbers and there is no indication of singularity or excessive pivot growth. If, however, underflows are handled in the usual way and set to zero, then no such simple guarantee can be made, and some kind of testing on the scaling of the problem is necessary. These results demonstrate that gradual underflow makes it easier to write reliable linear equation solvers than "store zero."

Appendix 1: Program Listing of Triangular Decomposition

```
array A[1..N,1..N] of real;
array L[1..N,1..N] of real;
array U[1..N,1..N] of real;
/* note that arrays A,L, and U can occupy the same
locations in memory; we use three different arrays
here to keep the notation consistent with
the previous analysis */
```

```
/* Triangular Factorization: A=LU */
/* As in the analysis, we assume the matrix has
been permuted so row or column exchanges are
unnecessary. */
/* Lines prefixed by "G.U.:" are the additions to
the code required for reliability using gradual
underflow. */
```

```
/*set mode to prevent spurious warnings about
denormalized numbers */
G.U.: invoke normalizing mode
```

```
for p=1 to N do /* p=pivot number */
begin
for j=p to N do /* compute U[p,j] */
begin
SUM=A[p,j];
for k=1 to p-1 do SUM=SUM - L[p,k]*U[k,j];
U[p,j]=SUM;
end
/* insert your favorite test for singularity
or pivot growth */
for i=p+1 to N do /* compute L[i,p] */
begin
SUM=A[i,p];
for k=1 to p-1 do SUM=SUM - L[i,k]*U[k,p];
L[i,p]=SUM/U[p,p];
end
end
```

Appendix 2: Program Listing of Solution of Both Triangular Systems Following Triangular Decomposition

```
array B[1..N] of real;
array X[1..N] of real;
array Y[1..N] of real;
/* Again, B,X, and Y could occupy the same locations
in memory. Arrays A,L and U are defined above. */
```

```

/* Forward Substitution: LY=B */
for i=1 to N do
  begin
  SUM=0;
  for k=1 to i-1 do SUM=SUM + L[i,k]*Y[k];
  Y[i]=B[i]-SUM;
  end

/* Back Substitution: UX=Y */

/* initialize flag for reporting underflows
in the solution x */
G.U.: xunderflow = false

for i=N downto 1 do
  begin
  SUM=0;
  for k=i+1 to N do SUM=SUM + U[i,k]*X[k];
  temp=Y[i] - SUM
  G.U.:underflowflag = false;
  X[i]=temp/U[i,i];
  G.U.:if underflowflag then xunderflow = underflowflag;
  end

/* if xunderflow = true, an underflow has
occurred in the solution X */

```

References

- ¹J. T. Coonen, "Underflow and the Denormalized Numbers", Computer, vol. 14, no. 3, March 1981, pp 75-87
- ²J. J. Dongarra, J. R. Bunch, C. B. Moler, G. W. Stewart, LINPACK User' Guide, Society for Industrial and Applied Mathematics, Philadelphia, 1979
- ³W. Kahan, J. Palmer, "On a Proposed Floating-Point Standard", SIGNUM Newsletter, October 1979
- ⁴J. Coonen, et al., "A Proposed Standard for Binary Floating Point Arithmetic", SIGNUM Newsletter, October 1979
- ⁵R. D. Skeel, "Scaling for Numerical Stability in Gaussian Elimination", Journal of the ACM, vol. 28, no. 3, July 1979
- ⁶VAX is a trademark of the Digital Equipment Corporation.
- ⁷J.H. Wilkinson, "Rounding Errors in Algebraic Processes", Prentice Hall, 1963
- ⁸W. Oettli and W. Prager, "Compatibility of Approximate Solution of Linear Equations with Given Error Bounds for Coefficients and Right-hand Sides," Numer. Math. 6(1964), 405-409
- ⁹J. W. Demmel, "Effects of Underflow on Solving Linear Systems", Computer Science Division, Electronic Research Laboratory Internal Report, U.C. Berkeley, 1981

Acknowledgement

The author gratefully acknowledges the support of the U. S. Department of Energy (Contract DE-AM03-76SF00034, Project Agreement DE-AS03-79ER10358), the Office of Naval Research (Contract N00014-76-C-0013), and Prof. W. Kahan, whose comment and encouragement have been invaluable.