

ON A MODEL FOR COMPUTING ROUND-OFF ERROR OF A SUM



AD 713972

BY  
GEORGE B. DANTZIG

*Handwritten initials*

STAN-CS-70-156  
MARCH 1970

Reproduced by  
NATIONAL TECHNICAL  
INFORMATION SERVICE  
Springfield, Va. 22151

DDC  
RECEIVED  
NOV 6 1970  
RECEIVED  
B

*Handwritten initials*

COMPUTER SCIENCE DEPARTMENT  
School of Humanities and Sciences  
STANFORD UNIVERSITY



**ON A MODEL FOR COMPUTING ROUND-OFF ERROR OF A SUM**

**BY**

**GEORGE B. DANTZIG**

**March, 1970**

**Technical Report No. 156**

**Computer Science Department  
Stanford University  
Stanford, California**

Reproduction and research of this report was partially supported by Office of Naval Research, Contract N00014-67-A0112-0011; U.S. Atomic Energy Commission, Contract AT[04-3] 326 PA #18; National Science Foundation, Grant GP 9329 and Grant GJ 320; U.S. Army Research Office, Contract DAHC 04-67-0028; and National Institutes of Health, Grant GM 14789-02.

Reproduction in whole or in part is permitted for any purpose of the United States Government. This document has been approved for public release and sale; its distribution is unlimited.

→ a sub-1, a sub-2, ..., a sub-n, the report discuss

→ Given real numbers  $a_1, a_2, \dots, a_n$ , we are interested in the classic problem of the error in computing  $S = \sum_{i=1}^n a_i$  when the sum is computed by  $\tilde{S}_0 = \sum_{i=1}^n a_i^*$  where  $a_i^*$  is the nearest integer to  $a_i$ . We shall first study this error as a function of a  $\Delta$  shift, i.e., when all numbers  $a_i$  are each shifted  $\Delta$  and then rounded;

$$(1) \quad S - n\Delta = \sum_{i=1}^n (a_i - \Delta)$$

$$(2) \quad \tilde{S}_\Delta - n\Delta = \sum_{i=1}^n (a_i - \Delta)^*$$

We will then let  $\Delta$  become a random variable that can take on uniformly any value in the interval  $-\frac{1}{2} \leq \Delta \leq +\frac{1}{2}$ . Different choices of  $\Delta$  give rise to different rounding errors  $\tilde{S}_\Delta - S$  and the variance of the distribution of  $\tilde{S}_\Delta - S$  can be used to measure the variability of the rounding error due to the random selection of the origin of the real numbers  $a_i$  with respect to that of the computer.

The cumulative error from (1) and (2) is

$$(3) \quad \tilde{S}_\Delta - S = \sum_{i=1}^n [(a_i - \Delta)^* - (a_i - \Delta)]$$

Let  $f_i$  be the positive fractional part of  $a_i$  and let  $\alpha_i$  be the largest integer not exceeding  $a_i$ , i.e.,

$$(4) \quad a_i = \alpha_i + f_i$$

Denoting by  $r_i$  the error of the  $i^{\text{th}}$  term, we have

$$(5) \quad r_i = [(a_i - \Delta)^* - (a_i - \Delta)] = \begin{cases} 1 - (f_i - \Delta) & \text{if } -\frac{1}{2} \leq \Delta \leq -\frac{1}{2} + f_i \\ -(f_i - \Delta) & \text{if } -\frac{1}{2} + f_i \leq \Delta \leq +\frac{1}{2} \end{cases}$$

To prove the above, we note that  $f_i - \Delta = (a_i - \Delta) + \alpha_i$ . If  $-\frac{1}{2} \leq f_i - \Delta \leq +\frac{1}{2}$  then  $(a_i - \Delta)$  is rounded to  $\alpha_i$ . Hence  $a_i - \Delta$  is rounded down if  $-\frac{1}{2} + f_i \leq \Delta$  otherwise rounded up.

Denoting expected value by  $E$ , we have by direct evaluation

$$(6) \quad E(r_i) = \int_{-\frac{1}{2}}^{+\frac{1}{2}} r_i d\Delta = 0$$

Assume  $f_i \leq f_j$ , then

$$\begin{aligned} E(r_i r_j) &= \int_{-\frac{1}{2}}^{-\frac{1}{2}+f_i} r_i r_j d\Delta + \int_{-\frac{1}{2}+f_i}^{-\frac{1}{2}+f_j} r_i r_j d\Delta + \int_{-\frac{1}{2}+f_j}^{+\frac{1}{2}} r_i r_j d\Delta \\ &= \int_{-\frac{1}{2}}^{+\frac{1}{2}} (f_i f_j - \Delta(f_i + f_j) + \Delta^2) d\Delta \\ &\quad + \int_{-\frac{1}{2}}^{-\frac{1}{2}+f_i} [(1-f_i-f_j) + 2\Delta] d\Delta \\ &\quad + \int_{-\frac{1}{2}+f_i}^{-\frac{1}{2}+f_j} [-f_i + \Delta] d\Delta \end{aligned}$$

Performing indicated integration yields:

$$(7) \quad E(r_i r_j) = \frac{1}{2}[|f_j - f_i|^2 - |f_j - f_i| + \frac{1}{6}]$$

which is one-half the 2<sup>nd</sup> order Bernoulli Polynomial in  $|f_j - f_1|$ . For  $f_j < f_1$  we also get (7). Note that the individual errors  $r_1$  and  $r_j$  are not independent of one another.

It now follows that

$$(9) \quad E(\tilde{S}) = S$$

$$(10) \quad E(\tilde{S}-S)^2 = E\left(\sum_{i=1}^n \sum_{j=1}^n r_i r_j\right) = \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n [ |f_1 - f_j|^2 - |f_1 - f_j| + \frac{1}{6} ]$$

The usual value of variance,  $E(\tilde{S}-S)^2 = n/12$ , will result if we further assume  $f_1$  are independently drawn from uniform distributions on  $[0 \leq f_1 \leq 1]$ .

Theorem: If the fractional parts of all  $a_1$  are equal to each other, then each term of (10) is maximum for  $0 \leq f_1 \leq 1$  and

$$(11) \quad \text{Max } E(S-S)^2 = \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \left(\frac{1}{6}\right) = \frac{n^2}{12}$$

From (10) we have an interesting inequality, namely for all  $f_1$

$$(12) \quad V(f) = \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \{ |f_1 - f_j|^2 - |f_1 - f_j| + \frac{1}{6} \} \geq 0$$

This function is not convex even for  $n=2$ , since  $f^{(1)} = (\frac{1}{2}, 0)$  and  $f^{(2)} = (-\frac{1}{2}, 0)$  yields  $V(f^0) = V(f^1) = \frac{1}{12} + \frac{1}{12} - \frac{1}{12} = \frac{1}{12}$  but

$V\left(\frac{f^1 + f^2}{2}\right) = V(0) = \frac{3}{12}$ . There appears to be no obvious direct way to establish that  $V(f) \geq 0$  for all  $0 \leq f_1 \leq 1$ . Our development shows  $V(f)$  to be a variance and this, of course, constitutes an indirect proof. We can replace (12) by a convex realization: Assume  $f_1 \geq f_{i+1}$  for all  $i$ , then the problem of finding  $\text{Min } V(f)$  can be rewritten:

$$\begin{aligned}
 (13) \quad \text{Find Min } [V(f)] = & \sum_{i < j} (f_i - f_j)^2 + \frac{n^2}{12} - [(n-1)f_1 + (n-3)f_2 + (n-5)f_3 \\
 & + \dots + (n-2k+1)f_k + \dots - (n-1)f_n]
 \end{aligned}$$

subject to

$$(14) \quad f_1 \geq f_2 \dots \geq f_n$$

$$(15) \quad 0 \leq f_1 \leq 1$$

Formally (13), (14), (15), is a positive definite quadratic program. Fortunately, as we shall see this can be solved by classical calculus by ignoring inequalities (14) and (15).

Theorem: Equally spaced  $f_i = (i - 1)/n$ , ( $i = 1, \dots, n$ ) yields  $\text{Min } V(f) = \frac{1}{12}$  independent of  $n$ , i.e., the variance of the sum in this case is minimum and is the same as the variance of the individual terms forming the sum.

**Proof:** Setting partials = 0 in (13) yields:

$$(16) \left\{ \begin{array}{l} 2(n-1)f_1 - 2f_2 \dots - 2f_n = (n-1) \\ - 2f_1 + 2(n-1)f_2 - 2f_n = (n-3) \\ - 2f_1 - 2f_2 \dots 2(n-1)f_{n-1} - 2f_n = -(n-3) \\ - 2f_1 - 2f_2 \dots 2(n-1)f_n = -(n-1) \end{array} \right.$$

Adding shows the equations to be dependent. Hence we may drop the last equation as redundant. Moreover, we can always translate the  $f_i$  so that the smallest  $f_i$ , namely  $f_n = 0$

Re-adding yields:

$$2f_1 + 2f_2 + \dots + 2f_{n-1} + 0 = (n-1), \quad f_n = 0.$$

Adding this last equation to each of the others gives

$$2nf_1 = (n - 2i + 1) + (n - 1) = 2(n - i)$$

$$(17) \quad f_1 = (n - i)/n$$

Evidently the conditions  $0 \leq f_i \leq 1$  and  $f_i \geq f_{i+1}$  are (by good luck) also satisfied so that (17) yields the minimum, namely

$$(18) \quad \text{Min } V(f) = \frac{n^2}{12} - \frac{1}{2} \sum_{i=1}^n (n-2_i+1)f_i = \frac{1}{12} .$$