

Assignment 1.

Given January 24, due January 30.

Objective: To explore computer arithmetic.

1. The fibonacci numbers, f_k , are defined by $f_0 = 1$, $f_1 = 1$, and

$$f_{k+1} = f_k + f_{k-1} \quad (1)$$

for any integer $k > 1$. A small perturbation of them, the “pib numbers” (“p” instead of “f” to indicate either the pentium bug or perturbation), p_k , are defined by $p_0 = 1$, $p_1 = 1$, and

$$p_{k+1} = c \cdot p_k + p_{k-1} \quad (2)$$

for any integer $k > 1$, where $c = 1 + \sqrt{3}/100$.

- a. Make a plot of $\log(f_n)$ and $\log(p_n)$ as a function of n . On the plot, mark $1/\epsilon_{mach}$ for single and double precision IEEE floating point arithmetic. This can be useful in answering the questions below.
 - b. For various n values, compute the f_k for $k = 2, 3, \dots, n$ using (1). Then rewrite (1) to express f_{k-1} in terms of f_k and f_{k+1} . Use the computed f_n and f_{n-1} to recompute f_k for $k = n-2, n-3, \dots, 0$. Make a plot of the difference between the original $f_0 = 1$ and the recomputed f_0 as a function of n . What n values result in no accuracy for the recomputed f_0 ? How do the results in single and double precision differ?
 - c. Repeat b. for the pib numbers. Comment on the striking difference in the way precision is lost in these two cases. Which is more typical? *Extra credit:* predict the order of magnitude of the error in recomputing p_0 using what you may know about recurrence relations and what you should know about computer arithmetic.
2. The binomial coefficients, $a_{n,k}$, are defined by

$$a_{n,k} = \binom{n}{k} = \frac{n!}{k!(n-k)!} \quad (3)$$

To compute the $a_{n,k}$, for a given n , start with $a_{n,0} = 1$ and then use the recurrence relation $a_{n,k+1} = \frac{n-k}{k+1}a_{n,k}$.

- (a) For a fixed of n , compute the $a_{n,k}$ this way, noting the largest $a_{n,k}$ and the accuracy with which $a_{n,n} = 1$ is computed. Do this in single and double precision. Why is it that roundoff is not a problem here as it was in problem (1)?

- b.** Use the algorithm of part (a) to compute

$$E(k) = \frac{1}{2^n} \sum_{k=0}^n k a_{n,k} = \frac{n}{2} . \quad (4)$$

In this equation, we think of k as a random variable, the number of heads in n tosses of a fair coin, with $E(k)$ being the expected value of k . This depends on n . Write a program without any safeguards against overflow or zero divide (*this time only!*)¹. Show (both in single and double precision) that the computed answer has high accuracy as long as the intermediate results are within the range of floating point numbers. As with (a), explain how the computer gets an accurate, small, answer when the intermediate numbers have such a wide range of values. Why is cancellation not a problem? Note the advantage of a wider range of values: we can compute $E(k)$ for much larger n in double precision. Print $E(k)$ as computed by (4) and $M_n = \max_k a_{n,k}$. For large n , one should be `inf` and the other `NaN`. Why?

- c.** For fairly large n , plot $a_{n,k}$ as a function of k for a range of k chosen to illuminate the interesting “bell shaped” behavior of the $a_{n,k}$ near their maximum.
- d.** (*Extra credit, and lots of it!*) Find a way to compute

$$S(k) = \sum_{k=0}^n (-1)^k \sin(2\pi \sin(k/n)) a_{n,k}$$

with good relative accuracy for large n . This is very hard, so don't spend too much time on it.

- 3.** In the example from Section 3 of the *Computer Arithmetic* chapter, show that the problem of computing x_1^0 and x_2^0 from x_1^k and x_2^k is ill conditioned for large k , under the condition that x_1^0 and x_2^0 are both positive and within a factor of 2 of each other. Identify what all the terms on the right of (5) are in this case. By contrast, show that computing y_1^0 and y_2^0 from y_1^k and y_2^k is well conditioned even for large k . This shows that if you actually knew y_1^k and y_2^k to high relative accuracy, you would be able to reconstruct x^0 to high relative accuracy. Unfortunately, computing y^k to high relative accuracy from x^k (with small relative errors) is impossible.

¹One of the purposes of the IEEE floating point standard was to *allow* a program with overflow or zero divide to run and print results