

---

# Numerical Analysis

Math 370 Spring 2009

MWF 11:30am - 12:25pm Fowler 110

©2009 Ron Buckmire

<http://faculty.oxy.edu/ron/math/370/09/>

---

## *Class 2*

**SUMMARY** (Un)Avoidable Errors in Computing With Machines

**CURRENT READING** Mathews & Fink, Sec 1.3

### **WARM-UP**

Write down the meaning of the terms **overflow**, **underflow**, **mantissa**, **characteristic**, **hexadecimal**, **bit**, **decimal machine number**, **(numerical) precision** and compare your answers with at least one other person in the class (after you have written your own!)

### **RECALL**

***k*-th digit Chopping**

In this case all the digits after  $d_k$  are **ignored** (“chopped off”)

***k*-th digit Rounding**

In this case if the value of  $d_{k+1} \geq 5$  then  $d_k$  is replaced by  $d_k + 1$

**Absolute Error**

If  $\tilde{p}$  is an approximation to  $p$ , the **absolute error** is  $|\tilde{p} - p|$

**Relative Error**

Provided  $p \neq 0$ , the **relative error** is  $\frac{|\tilde{p} - p|}{|p|}$

### **EXAMPLE**

Show that the expression involving  $k$  which gives you an upper bound for the relative error involved in using chopping arithmetic is  $\epsilon_{rel} = 10^{-k+1}$

It can also be shown that a bound for the relative error involved in using **rounding arithmetic** is *half* that for chopping,  $\epsilon_{rel} = 0.5 \times 10^{-k+1} = 5 \times 10^{-k}$ .

**Significant Digits**

The number  $\tilde{p}$  is said to approximate  $p$  to  $k$  significant digits (or figures) if  $k$  is the largest non-negative integer for which the relative error is less than  $5 \times 10^{-k}$ . (Mathews p. 25)

## Round-off Errors in the Quadratic Formula

Recall that the common formula for the roots of a quadratic equation  $ax^2 + bx + c = 0$  is

$$x_1 = \frac{-b + \sqrt{b^2 - 4ac}}{2a} \quad \text{and} \quad x_2 = \frac{-b - \sqrt{b^2 - 4ac}}{2a}$$

Round-off error can wreak havoc with the numerical implementation of this formula. Consider

$$x^2 + 62.10x + 1 = 0$$

which has the approximate roots  $x_1 = -0.01610723$  and  $x_2 = -62.08390$

Because of the size of the parameters in the quadratic equation,  $b^2$  is much bigger than  $4ac$ , so  $\sqrt{b^2 - 4ac}$  is very close to  $b$ .

$$a = 1, b = 62.10, c = 1$$

$$b^2 = \qquad 4ac = \qquad b^2 - 4ac = \qquad \sqrt{b^2 - 4ac} =$$

### GROUPWORK

Using 4-digit rounding (or chopping) arithmetic compute the first root  $x_1$

What's the relative error in this calculation?

Solution: change the formula for  $x_1$  so that we don't have to subtract  $b$  from  $\sqrt{b^2 - 4ac}$

Now, a new formula for  $x_1 =$

Use a similar new formula to compute  $x_2$  (using 4-digit precision) and compute the relative error in  $x_2$

What's the problem?

Solution: Use the new formula for  $x_1$  when you have to subtract numbers which are similar in size, use the traditional formula for the other root.

When you subtract numbers of very similar size there will be a **loss of significance** or **subtractive cancellation** of digits when this operation is computed. **AVOID DOING THIS!**

## The Ultimate Quadratic Formula

$$q \equiv -\frac{1}{2} [b + \text{sign}(b)\sqrt{b^2 - 4ac}]$$

where

$$\text{sign}(b) = \begin{cases} 1 & b \geq 0 \\ -1 & b < 0 \end{cases}$$

and

$$x_1 = \frac{q}{a} \quad \text{and} \quad x_2 = \frac{c}{q}$$