

# Numerical Errors

- Introduction
- Representation of Numbers in a Computer
- Numerical Differentiation
- Factorable Functions and Algorithmic Differentiation

# Contents

- Introduction
- Representation of Numbers in a Computer
- Numerical Differentiation
- Factorable Functions and Algorithmic Differentiation

# Objectives

In this lecture we will learn about

- the fact that many computer programs store numbers with finite precision only
- the IEEE standard for storing floating point numbers
- numerical and algorithmic differentiation

# Contents

- Introduction
- Representation of Numbers in a Computer
- Numerical Differentiation
- Factorable Functions and Algorithmic Differentiation

# Scientific Computing

Computers or calculators typically store numbers with finite precision:

- Example 1:  $8 + 8 == 16$  ?
- Example 2:  $(\sqrt{5})^2 == 5$  ?
- Example 3:  $1.1 + 0.1 == 1.2$  ?

Let's try this with JULIA:

```
julia> (1.1 + 0.1) == 1.2  
false
```

```
julia> 1.1 + 0.1  
1.2000000000000002
```

Problem: numerical error:  $\approx 2 * 10^{-16}$ .

# Scientific Computing

Computers or calculators typically store numbers with finite precision:

- Example 1:  $8 + 8 == 16$  ?
- Example 2:  $(\sqrt{5})^2 == 5$  ?
- Example 3:  $1.1 + 0.1 == 1.2$  ?

Let's try this with JULIA:

```
julia>(1.1 + 0.1) == 1.2  
false
```

```
julia> 1.1 + 0.1  
1.2000000000000002
```

Problem: numerical error:  $\approx 2 * 10^{-16}$ .

# Scientific Computing

Computers or calculators typically store numbers with finite precision:

- Example 1:  $8 + 8 == 16$  ?
- Example 2:  $(\sqrt{5})^2 == 5$  ?
- Example 3:  $1.1 + 0.1 == 1.2$  ?

Let's try this with JULIA:

<pre>julia&gt;(1.1 + 0.1) == 1.2 false</pre>	<pre>julia&gt; 1.1 + 0.1 1.2000000000000002</pre>
--	---

Problem: numerical error:  $\approx 2 * 10^{-16}$ .

# Floating Point Numbers

IEEE standard for double-precision floating point numbers:

$$x = \pm(1 + m) \cdot 2^e \quad \text{with} \quad m = \sum_{i=1}^{52} m_i 2^{-i} \quad \text{and} \quad e = \sum_{i=0}^{10} c_i 2^i - \bar{c},$$

Names:  $m$  = mantissa,  $e$  = exponent.

Storage requirement:

- 1 bit to store the sign.
- 11 bits to store  $c_{10}, \dots, c_0 \in \{0, 1\}$ ; offset  $\bar{c} = 1023$ .
- 52 bits to store  $m_1, \dots, m_{52} \in \{0, 1\}$  .

In total:  $(1 + 11 + 52)$  bits = 64 bits = 8 bytes.



# Floating Point Numbers

IEEE standard for double-precision floating point numbers:

$$x = \pm(1 + m) \cdot 2^e \quad \text{with} \quad m = \sum_{i=1}^{52} m_i 2^{-i} \quad \text{and} \quad e = \sum_{i=0}^{10} c_i 2^i - \bar{c},$$

Names:  $m$  = mantissa,  $e$  = exponent.

Storage requirement:

- 1 bit to store the sign.
- 11 bits to store  $c_{10}, \dots, c_0 \in \{0, 1\}$ ; offset  $\bar{c} = 1023$ .
- 52 bits to store  $m_1, \dots, m_{52} \in \{0, 1\}$  .

In total:  $(1 + 11 + 52)$  bits = 64 bits = 8 bytes.

# Floating Point Numbers

IEEE standard for double-precision floating point numbers:

$$x = \pm(1 + m) \cdot 2^e \quad \text{with} \quad m = \sum_{i=1}^{52} m_i 2^{-i} \quad \text{and} \quad e = \sum_{i=0}^{10} c_i 2^i - \bar{c},$$

Names:  $m$  = mantissa,  $e$  = exponent.

Storage requirement:

- 1 bit to store the sign.
- 11 bits to store  $c_{10}, \dots, c_0 \in \{0, 1\}$ ; offset  $\bar{c} = 1023$ .
- 52 bits to store  $m_1, \dots, m_{52} \in \{0, 1\}$  .

In total:  $(1 + 11 + 52)$  bits = 64 bits = 8 bytes.

## A closer look at the example

```
julia>bits(1.1)
```

```
"0011111111110001100110011001100110011001100110011001100110011010"
```

```
julia>bits(0.1)
```

```
"001111111011100110011001100110011001100110011001100110011010"
```

```
julia>bits(1.2)
```

```
"001111111111001100110011001100110011001100110011001100110011"
```

```
julia>bits(1.1+0.1)
```

```
"001111111111001100110011001100110011001100110011001100110100"
```

## A closer look at the example

```
julia>bits(1.1)
```

```
"001111111111000110011001100110011001100110011001100110011010"
```

```
julia>bits(0.1)
```

```
"001111111011100110011001100110011001100110011001100110011010"
```

```
julia>bits(1.2)
```

```
"001111111111001100110011001100110011001100110011001100110011"
```

```
julia>bits(1.1+0.1)
```

```
"001111111111001100110011001100110011001100110011001100110100"
```

## A closer look at the example

```
julia>bits(1.1)
```

```
"001111111110001100110011001100110011001100110011001100110011010"
```

```
julia>bits(0.1)
```

```
"001111111011100110011001100110011001100110011001100110011010"
```

```
julia>bits(1.2)
```

```
"00111111111001100110011001100110011001100110011001100110011"
```

```
julia>bits(1.1+0.1)
```

```
"001111111111001100110011001100110011001100110011001100110100"
```

## A closer look at the example

```
julia>bits(1.1)
```

```
"0011111111110001100110011001100110011001100110011001100110011010"
```

```
julia>bits(0.1)
```

```
"001111111011100110011001100110011001100110011001100110011010"
```

```
julia>bits(1.2)
```

```
"001111111111001100110011001100110011001100110011001100110011"
```

```
julia>bits(1.1+0.1)
```

```
"001111111111001100110011001100110011001100110011001100110100"
```

# Floating Point Numbers

- Numbers between 1 and  $1 + 2^{-52}$  cannot be represented.
- The (relative) rounding  $\text{eps} = 2^{-52}$  is called *machine precision*.
- The absolute rounding error  $\text{eps} * 2^e$  depends on exponent  $e$ .  
(if we work with larger numbers, we get larger rounding errors)

**Important to remember:**  $\text{eps} = 2^{-52} \approx 2 * 10^{-16}$ .

# Floating Point Numbers

- Numbers between 1 and  $1 + 2^{-52}$  cannot be represented.
- The (relative) rounding  $\text{eps} = 2^{-52}$  is called *machine precision*.
- The absolute rounding error  $\text{eps} * 2^e$  depends on exponent  $e$ .  
(if we work with larger numbers, we get larger rounding errors)

**Important to remember:**  $\text{eps} = 2^{-52} \approx 2 * 10^{-16}$ .



# Numerical function evaluation

Let us evaluate the function

$$\Phi(x) = \sin(10^8 x)$$

at  $x = \pi$ . The exact solution is  $\Phi(\pi) = 0$ .

```
julia> sin(10^8 pi)  
-3.9082928156687315e - 8
```

**Caution:** The function  $\Phi$  is ill-conditioned, i.e., the evaluation error is much larger than  $\text{eps} \approx 2 * 10^{-16}$ !

# Numerical function evaluation

Let us evaluate the function

$$\Phi(x) = \sin(10^8 x)$$

at  $x = \pi$ . The exact solution is  $\Phi(\pi) = 0$ .

```
julia> sin(10^8 pi)  
-3.9082928156687315e - 8
```

**Caution:** The function  $\Phi$  is ill-conditioned, i.e., the evaluation error is much larger than  $\text{eps} \approx 2 * 10^{-16}$ !

# Numerical function evaluation

Let us evaluate the function

$$\Phi(x) = \sin(10^8 x)$$

at  $x = \pi$ . The exact solution is  $\Phi(\pi) = 0$ .

```
julia> sin(10^8 pi)  
-3.9082928156687315e - 8
```

**Caution:** The function  $\Phi$  is ill-conditioned, i.e., the evaluation error is much larger than  $\text{eps} \approx 2 * 10^{-16}$ !

# Storing Numbers in a Computer

There exists a variety of ways to represent numbers:

- Floating point numbers be it 64bit (“double precision”) or 32bit (“single precision”).
- Integers are often stored differently. Remark:  
`julia>bits(3)` is not the same as `julia>bits(3.)!!!`
- Arbitrary precision arithmetics are an alternative (not our focus).
- Verified arithmetics store intervals rather than single numbers.

# Storing Numbers in a Computer

There exists a variety of ways to represent numbers:

- Floating point numbers be it 64bit (“double precision”) or 32bit (“single precision”).

- Integers are often stored differently. Remark:

`julia>bits(3)` is not the same as `julia>bits(3.)!!!`

- Arbitrary precision arithmetics are an alternative (not our focus).
- Verified arithmetics store intervals rather than single numbers.

# Storing Numbers in a Computer

There exists a variety of ways to represent numbers:

- Floating point numbers be it 64bit (“double precision”) or 32bit (“single precision”).
- Integers are often stored differently. Remark:  
`julia>bits(3)` is not the same as `julia>bits(3.)!!!`
- Arbitrary precision arithmetics are an alternative (not our focus).
- Verified arithmetics store intervals rather than single numbers.

# Storing Numbers in a Computer

There exists a variety of ways to represent numbers:

- Floating point numbers be it 64bit (“double precision”) or 32bit (“single precision”).
- Integers are often stored differently. Remark:  
`julia>bits(3)` is not the same as `julia>bits(3.)`!!!
- Arbitrary precision arithmetics are an alternative (not our focus).
- Verified arithmetics store intervals rather than single numbers.

# Contents

- Introduction
- Representation of Numbers in a Computer
- Numerical Differentiation
- Factorable Functions and Algorithmic Differentiation



# Finite Differences

The derivative of a twice continuously differentiable function  $f : \mathbb{R} \rightarrow \mathbb{R}$  can be approximated by finite differences:

$$\frac{\partial f(x)}{\partial x} \approx \frac{f(x+h) - f(x)}{h}$$

- The mathematical approximation error, given by

$$\left| \frac{f(x+h) - f(x)}{h} - \frac{\partial f}{\partial x}(x) \right| \approx \frac{h}{2} \left| \frac{\partial^2 f}{\partial x^2}(x) \right| = \mathbf{O}(h),$$

tends to zero for  $h \rightarrow 0$ .

- The numerical error is approximately

$$\mathbf{O}\left(\frac{\text{eps}}{h}\right)$$

- If  $f$  is well conditioned, we choose

$$h \approx \operatorname{argmin}_h \left( h + \frac{\text{eps}}{h} \right) = \sqrt{\text{eps}}.$$

# Finite Differences

The derivative of a twice continuously differentiable function  $f : \mathbb{R} \rightarrow \mathbb{R}$  can be approximated by finite differences:

$$\frac{\partial f(x)}{\partial x} \approx \frac{f(x+h) - f(x)}{h}$$

- The mathematical approximation error, given by

$$\left| \frac{f(x+h) - f(x)}{h} - \frac{\partial f}{\partial x}(x) \right| \approx \frac{h}{2} \left| \frac{\partial^2 f}{\partial x^2}(x) \right| = \mathbf{O}(h),$$

tends to zero for  $h \rightarrow 0$ .

- The numerical error is approximately

$$\mathbf{O}\left(\frac{\text{eps}}{h}\right)$$

- If  $f$  is well conditioned, we choose

$$h \approx \operatorname{argmin}_h \left( h + \frac{\text{eps}}{h} \right) = \sqrt{\text{eps}}.$$

# Finite Differences

The derivative of a twice continuously differentiable function  $f : \mathbb{R} \rightarrow \mathbb{R}$  can be approximated by finite differences:

$$\frac{\partial f(x)}{\partial x} \approx \frac{f(x+h) - f(x)}{h}$$

- The mathematical approximation error, given by

$$\left| \frac{f(x+h) - f(x)}{h} - \frac{\partial f}{\partial x}(x) \right| \approx \frac{h}{2} \left| \frac{\partial^2 f}{\partial x^2}(x) \right| = \mathbf{O}(h),$$

tends to zero for  $h \rightarrow 0$ .

- The numerical error is approximately

$$\mathbf{O}\left(\frac{\text{eps}}{h}\right)$$

- If  $f$  is well conditioned, we choose

$$h \approx \operatorname{argmin}_h \left( h + \frac{\text{eps}}{h} \right) = \sqrt{\text{eps}}.$$

# Finite Differences

The derivative of a twice continuously differentiable function  $f : \mathbb{R} \rightarrow \mathbb{R}$  can be approximated by finite differences:

$$\frac{\partial f(x)}{\partial x} \approx \frac{f(x+h) - f(x)}{h}$$

- The mathematical approximation error, given by

$$\left| \frac{f(x+h) - f(x)}{h} - \frac{\partial f}{\partial x}(x) \right| \approx \frac{h}{2} \left| \frac{\partial^2 f}{\partial x^2}(x) \right| = \mathbf{O}(h),$$

tends to zero for  $h \rightarrow 0$ .

- The numerical error is approximately

$$\mathbf{O}\left(\frac{\text{eps}}{h}\right)$$

- If  $f$  is well conditioned, we choose

$$h \approx \operatorname{argmin}_h \left( h + \frac{\text{eps}}{h} \right) = \sqrt{\text{eps}}.$$

## Central Differences

In order to reduce the mathematical approximation error, we can use central differences

$$\frac{\partial f(x)}{\partial x} \approx \frac{f(x+h) - f(x-h)}{2h}$$

to approximate the derivative of  $f$ .

- The mathematical approximation error is now

$$\left| \frac{f(x+h) - f(x-h)}{2h} - \frac{\partial f}{\partial x}(x) \right| \leq \mathbf{O}(h^2).$$

- The numerical error is still in the order of

$$\frac{\text{eps}}{h} = \mathbf{O}\left(\frac{\text{eps}}{h}\right)$$

- In practice, if  $f$  is well conditioned, we choose  $h \approx \sqrt[3]{\text{eps}}$ .

## Central Differences

In order to reduce the mathematical approximation error, we can use central differences

$$\frac{\partial f(x)}{\partial x} \approx \frac{f(x+h) - f(x-h)}{2h}$$

to approximate the derivative of  $f$ .

- The mathematical approximation error is now

$$\left| \frac{f(x+h) - f(x-h)}{2h} - \frac{\partial f}{\partial x}(x) \right| \leq \mathbf{O}(h^2).$$

- The numerical error is still in the order of

$$\frac{\text{eps}}{h} = \mathbf{O}\left(\frac{\text{eps}}{h}\right)$$

- In practice, if  $f$  is well conditioned, we choose  $h \approx \sqrt[3]{\text{eps}}$ .

## Central Differences

In order to reduce the mathematical approximation error, we can use central differences

$$\frac{\partial f(x)}{\partial x} \approx \frac{f(x+h) - f(x-h)}{2h}$$

to approximate the derivative of  $f$ .

- The mathematical approximation error is now

$$\left| \frac{f(x+h) - f(x-h)}{2h} - \frac{\partial f}{\partial x}(x) \right| \leq \mathbf{O}(h^2).$$

- The numerical error is still in the order of

$$\frac{\text{eps}}{h} = \mathbf{O}\left(\frac{\text{eps}}{h}\right)$$

- In practice, if  $f$  is well conditioned, we choose  $h \approx \sqrt[3]{\text{eps}}$ .

## Central Differences

In order to reduce the mathematical approximation error, we can use central differences

$$\frac{\partial f(x)}{\partial x} \approx \frac{f(x+h) - f(x-h)}{2h}$$

to approximate the derivative of  $f$ .

- The mathematical approximation error is now

$$\left| \frac{f(x+h) - f(x-h)}{2h} - \frac{\partial f}{\partial x}(x) \right| \leq \mathbf{O}(h^2).$$

- The numerical error is still in the order of

$$\frac{\text{eps}}{h} = \mathbf{O}\left(\frac{\text{eps}}{h}\right)$$

- In practice, if  $f$  is well conditioned, we choose  $h \approx \sqrt[3]{\text{eps}}$ .



# Contents

- Introduction
- Representation of Numbers in a Computer
- Numerical Differentiation
- Factorable Functions and Algorithmic Differentiation

## Factorable Functions

Many (but not all) functions of our interest can be composed into a finite list of atom operations from a given library  $L$ , e.g.,  $L = \{+, -, *, \sin, \cos, \log, \dots\}$ .

### Example

- The function  $f(x) = \sin(x_1 * x_2) + \cos(x_1)$  will (internally) be evaluated as

$$a_1 = x_1 * x_2$$

$$a_2 = \sin(a_1)$$

$$a_3 = \cos(x_1)$$

$$a_4 = a_2 + a_3$$

$$f(x) = a_4 .$$

Here, the memory for  $a_1, \dots, a_4$  is (usually) allocated temporarily.

## Factorable Functions

Many (but not all) functions of our interest can be composed into a finite list of atom operations from a given library  $L$ ,

e.g.,  $L = \{+, -, *, \sin, \cos, \log, \dots\}$ .

### Example

- The function  $f(x) = \sin(x_1 * x_2) + \cos(x_1)$  will (internally) be evaluated as

$$a_1 = x_1 * x_2$$

$$a_2 = \sin(a_1)$$

$$a_3 = \cos(x_1)$$

$$a_4 = a_2 + a_3$$

$$f(x) = a_4 .$$

Here, the memory for  $a_1, \dots, a_4$  is (usually) allocated temporarily.

# Algorithmic Differentiation

In modern computer programs, algorithmic differentiation (AD) is used in order to avoid discretization errors. Let's try to understand the main idea of forward AD by looking at an example:

$$\begin{array}{lcl} a_0 & = & x \\ a_1 & = & a_0 * a_0 \\ a_2 & = & \sin(a_1) \\ a_3 & = & a_1 + a_2 \\ f(x) & = & a_3 . \end{array} \quad \begin{array}{lcl} b_0 & = & 1 \\ b_1 & = & a_0 * b_0 + b_0 * a_0 \\ b_2 & = & \cos(a_1) * b_1 \\ b_3 & = & b_1 + b_2 \\ f'(x) & = & b_3 . \end{array}$$

In practice, this is usually implemented by operator overloading.

# Algorithmic Differentiation

In modern computer programs, algorithmic differentiation (AD) is used in order to avoid discretization errors. Let's try to understand the main idea of forward AD by looking at an example:

$$\begin{array}{lcl} a_0 & = & x \\ a_1 & = & a_0 * a_0 \\ a_2 & = & \sin(a_1) \\ a_3 & = & a_1 + a_2 \\ f(x) & = & a_3 . \end{array} \quad \left| \quad \begin{array}{lcl} b_0 & = & 1 \\ b_1 & = & a_0 * b_0 + b_0 * a_0 \\ b_2 & = & \cos(a_1) * b_1 \\ b_3 & = & b_1 + b_2 \\ f'(x) & = & b_3 . \end{array} \right.$$

In practice, this is usually implemented by operator overloading.

# Summary

- Programs often store numbers with finite precision only.
- IEEE double precision floating point numbers:  $\text{eps} \approx 2 * 10^{-16}$ .
- Numerical differentiation is in general less accurate than algorithmic differentiation (AD).