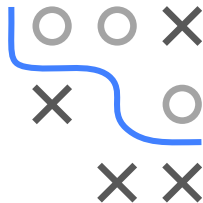


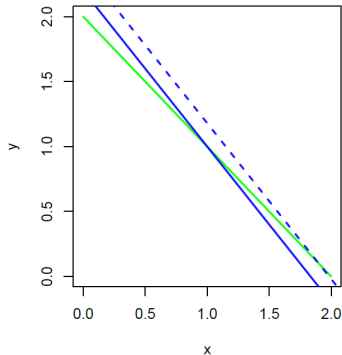
Algorithms and Data Structures

Numerics

Numerical Error & Conditioning



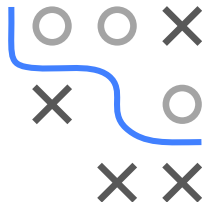
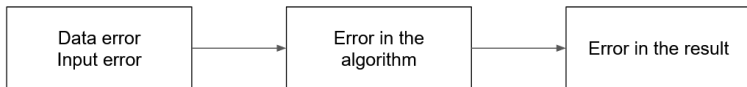
Ill-conditioned



Learning goals

- Error caused by condition of the problem
- Error caused by stability of algorithms

NUMERICAL ERROR



Errors in the result are caused by

- Data error or input error; often unavoidable, are part of the problem
→ **Condition of the problem**
- Error in the algorithm that can often be fixed by modification
→ **Stability of algorithms**

NUMERICAL ERROR / 2

Given: Approximate value \tilde{x} for exact value x .

- Absolute error (this is signed!):

$$\Delta x = \tilde{x} - x$$

If $|\Delta x| \leq \epsilon$, then ϵ is known as absolute error bound.

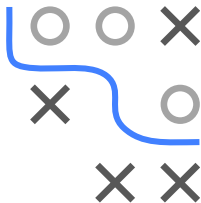
- Relative error:

$$\delta_x = \frac{|\tilde{x} - x|}{|x|} = \frac{|\Delta x|}{|x|}$$

If $\delta_x \leq \rho$, then ρ is known as relative error bound.

Comment:

Relative error is dimensionless, value in denominator may obviously not be too close to zero (\rightarrow use absolute error instead).

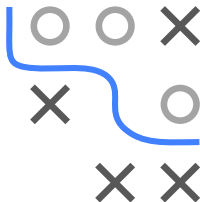


NUMERICAL ERROR / 3

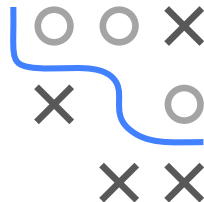
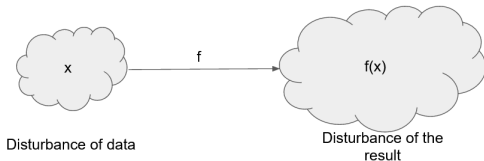
The following holds:

$$\tilde{x} = x \left(1 + \frac{\Delta x}{x} \right)$$

If only the first m digits of a number are determined in decimal representation, a possible relative error of 10^{-m} results (and vice versa).



CONDITION OF A PROBLEM



How sensitive is the result to a small disturbance of the data?

Well-conditioned: Result is a little sensitive

Ill-conditioned: Result is very sensitive

This is particularly important in statistics, since data is usually error-prone.

CONDITION OF A PROBLEM / 2

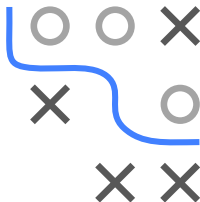
Example 1: System of equations

$$0.835x + 0.667y = 0.168$$

$$0.333x + 0.266y = 0.067$$

The solution is $x = 1$ and $y = -1$.

By a small change of the right side, for example $0.067 \rightarrow 0.066$, the solution changes to $x = -666$ and $y = 834$.



CONDITION FOR VECTORS

For vectors, errors are defined by means of a suitable norm.

Absolute error:

$$\Delta \mathbf{x} = \tilde{\mathbf{x}} - \mathbf{x}$$

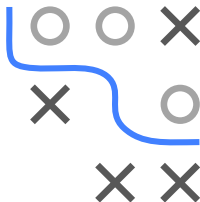
Relative error:

$$\delta = \frac{\|\tilde{\mathbf{x}} - \mathbf{x}\|}{\|\mathbf{x}\|}$$

Similarly, for $f : \mathbb{R}^n \rightarrow \mathbb{R}^d$ the **norm-wise** condition at the location \mathbf{x} is defined by using the condition number.

Condition number: The smallest $\kappa \geq 0$, such that

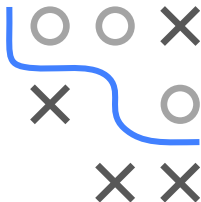
$$\frac{\|f(\mathbf{x} + \Delta \mathbf{x}) - \mathbf{y}\|}{\|\mathbf{y}\|} \leq \kappa \frac{\|\Delta \mathbf{x}\|}{\|\mathbf{x}\|} \text{ for } \|\Delta \mathbf{x}\| \rightarrow 0.$$



CONDITION FOR VECTORS / 2

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be differentiable in point x , the condition can be specified using the derivative:

$$\kappa = \frac{\|\mathbf{x}\|}{\|f(\mathbf{x})\|} \|\nabla f(\mathbf{x})\|.$$



Proof sketch:

The definition of the condition can be written as

$$\kappa = \limsup_{\|\Delta \mathbf{x}\| \rightarrow 0} \frac{\|f(\mathbf{x} + \Delta \mathbf{x}) - \mathbf{y}\|}{\|\mathbf{y}\|} \frac{\|\mathbf{x}\|}{\|\Delta \mathbf{x}\|} = \limsup_{\|\Delta \mathbf{x}\| \rightarrow 0} \frac{\|f(\mathbf{x} + \Delta \mathbf{x}) - \mathbf{y}\|}{\|\Delta \mathbf{x}\|} \frac{\|\mathbf{x}\|}{\|f(\mathbf{x})\|}$$

According to the mean value theorem the following applies: there is a $\mathbf{v} \in (\mathbf{0}, \Delta \mathbf{x})$ such that

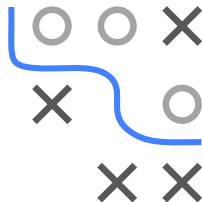
$$\nabla f(\mathbf{x} + \mathbf{v}) = \frac{\|f(\mathbf{x} + \Delta \mathbf{x}) - \mathbf{y}\|}{\|\Delta \mathbf{x}\|}.$$

CONDITION FOR VECTORS / 3

From $\|\Delta \mathbf{x}\| \rightarrow 0$ follows $\|\mathbf{v}\| \rightarrow 0$ and thus

$$\kappa = \limsup_{\|\Delta \mathbf{x}\| \rightarrow 0} \|\nabla f(\mathbf{x} + \mathbf{v})\| \frac{\|\mathbf{x}\|}{\|f(\mathbf{x})\|} = \|\nabla f(\mathbf{x})\| \frac{\|\mathbf{x}\|}{\|f(\mathbf{x})\|}$$

- $\nabla f(\mathbf{x})$ corresponds to the $1 \times n$ Jacobi matrix here. Thus, $\|\nabla f(\mathbf{x})\|$ is an induced *matrix norm*.
- The condition depends on the choice of the vector norm / induced matrix norm.



ARITHMETIC OPERATIONS

- **Multiplication:**

We look at $f(x_1, x_2) = x_1 \cdot x_2$ and calculate κ using

$$\kappa = \frac{\|\nabla f(\mathbf{x})\|_1 \cdot \|\mathbf{x}\|_1}{|f(\mathbf{x})|}$$

Since

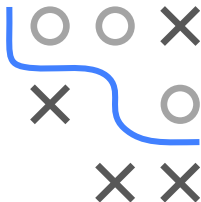
$$\|\nabla f(\mathbf{x})\|_1 = \|(x_2, x_1)\|_1 \stackrel{*}{=} \max(|x_1|, |x_2|)$$

$$\|\mathbf{x}\|_1 = \|(x_1, x_2)\|_1 = |x_1| + |x_2|$$

$$|f(\mathbf{x})| = |x_1 x_2|$$

we get $\kappa = 1 + \frac{\max(|x_1|, |x_2|)}{\min(|x_1|, |x_2|)}$. So κ becomes very big for $x_2 \ll x_1$.

* Attention: induced matrix norm (maximum absolute column sum norm).



ARITHMETIC OPERATIONS / 2

- **Addition / Subtraction:**

We will consider $f(x_1, x_2) = x_1 \pm x_2$ and calculate κ using

$$\kappa = \frac{\|\nabla f(\mathbf{x})\|_1 \cdot \|\mathbf{x}\|_1}{|f(\mathbf{x})|}$$

Since

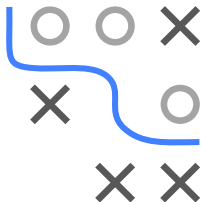
$$\|\nabla f(\mathbf{x})\|_1 = \|(1, \pm 1)\|_1^* = 1$$

$$\|\mathbf{x}\|_1 = \|(x_1, x_2)\|_1 = |x_1| + |x_2|$$

$$|f(\mathbf{x})| = |x_1 \pm x_2|$$

we get $\kappa = \frac{|x_1| + |x_2|}{|x_1 \pm x_2|}$.

* Attention: induced matrix norm (maximum absolute column sum norm).



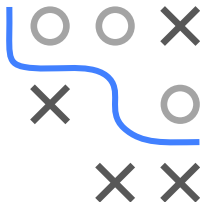
ARITHMETIC OPERATIONS / 3

The addition of two positive numbers is well-conditioned with

$$\kappa = \frac{|x_1| + |x_2|}{|x_1 + x_2|} = 1.$$

In the case of $x_1 \approx x_2$, the subtraction of two positive numbers is very ill-conditioned ($\kappa = \frac{|x_1| + |x_2|}{|x_1 + (-x_2)|}$ becomes arbitrarily large).

This problem is known as **loss of significance**.



LOSS OF SIGNIFICANCE (EXAMPLE) / 3

By subtracting the two numbers, the first 10 significant digits cancel each other out. The uncertainties from the 17th place now shift to the 7th significant place.

Due to the loss of significance we have lost $10 = 16 - 6$ digits (number of significant digits - number of deleted digits) in accuracy.

$$\begin{array}{r} 1.234567890123456???... \\ -1.234567890000000???... \\ \hline 0.000000000123456???... \\ = 1.23456???... \times 10^{-10} \end{array}$$

