

MATH3511/6111: Scientific Computing

09. Conditioning and Stability

Lindon Roberts

Semester 1, 2022

Office: Hanna Neumann Building #145, Room 4.87

Email: lindon.roberts@anu.edu.au

Based on lecture notes written by S. Roberts, L. Stals, Q. Jin, M. Hegland, K. Duru.



Australian
National
University

Conditioning and Stability

Throughout this course, we are looking at finding approximate solutions to mathematical problems.

In practice, we will then implement algorithms on a computer, which introduces rounding errors.

- **Conditioning** is about how sensitive **mathematical problems** are to perturbations.
- **Stability** is about how sensitive **algorithms** are to perturbations/rounding errors

Abstract Framework

All the mathematical processes we discuss can be viewed as a mapping from data to solution. We can formalise this:

- X is a vector space where our inputs (data) live.
- Y is a vector space where our solution lives.
- $f : X \rightarrow Y$ returns the true mathematical solution for given input data.

We are interested in how f behaves near our actual input data (call it \mathbf{x}).

Definition (Well-conditioned)

A problem is well-conditioned at \mathbf{x} if any small perturbations of the data \mathbf{x} cause only small changes in the solution $\mathbf{f}(\mathbf{x})$.

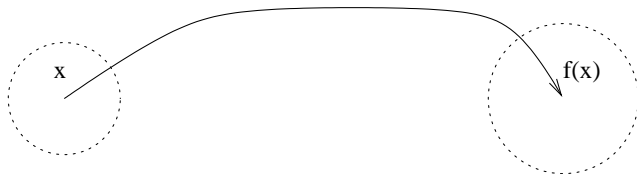


Figure 1. Well-conditioned problem: small perturbations in the data get mapped to small perturbations in the solution.

Conditioning

Definition (Ill-conditioned)

A problem is ill-conditioned at \mathbf{x} if **some** small perturbations of the data \mathbf{x} cause large changes in the solution $\mathbf{f}(\mathbf{x})$.

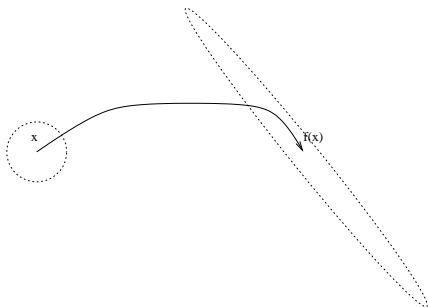


Figure 2. Ill-conditioned problem: some small perturbations in the data get mapped to large perturbations in the solution.

Condition Number

Question

How can we quantify the concept of well/ill-conditioned problems?

Suppose our original data \mathbf{x} is perturbed to $\mathbf{x} + \delta\mathbf{x}$. This gives us two (nearby) problems with solutions $\mathbf{f}(\mathbf{x})$ and $\mathbf{f}(\mathbf{x} + \delta\mathbf{x})$. The change in solution is

$$\delta\mathbf{f} = \mathbf{f}(\mathbf{x} + \delta\mathbf{x}) - \mathbf{f}(\mathbf{x}).$$

Definition (Absolute condition number)

The absolute condition number for the problem \mathbf{f} at a point \mathbf{x} is

$$\tilde{\kappa} = \max_{\delta\mathbf{x}} \frac{\|\delta\mathbf{f}\|}{\|\delta\mathbf{x}\|},$$

where the maximum is taken over all (infinitesimally small) $\delta\mathbf{x}$.

Recall from Taylor's theorem that if \mathbf{f} is differentiable,

$$\mathbf{f}(\mathbf{x} + \delta) = \mathbf{f}(\mathbf{x}) + J(\mathbf{x})\delta\mathbf{x} + \mathcal{O}(\|\delta\mathbf{x}\|^2),$$

where $J(\mathbf{x})$ is the Jacobian matrix of \mathbf{f} at \mathbf{x} .

Then since $\delta\mathbf{f} \approx J(\mathbf{x})\delta\mathbf{x}$, we have

$$\tilde{\kappa} = \|J(\mathbf{x})\|.$$

Condition Number

However, like all our other calculations, we always try to normalise our changes (e.g. relative errors). The relative condition number measure the relative changes in the solution vs. relative changes in the inputs.

Definition (Relative condition number)

The relative condition number for the problem \mathbf{f} at a point \mathbf{x} is

$$\kappa = \max_{\delta \mathbf{x}} \frac{\|\delta \mathbf{f}\| / \|\mathbf{f}(\mathbf{x})\|}{\|\delta \mathbf{x}\| / \|\mathbf{x}\|},$$

where the maximum is taken over all (infinitesimally small) $\delta \mathbf{x}$.

If \mathbf{f} is differentiable, then

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

Warning

The condition number only depends on the actual mathematical problem we care about. What algorithm we choose to find the solution is irrelevant!

It is a fundamental part of the problem: no algorithm can get around ill-conditioning (without changing the mathematical problem it solves).

In general, small κ means well-conditioned and large κ means ill-conditioned, but there is no clear cutoff between them.

Condition Number: Examples

Problem

Given a positive number, multiply it by 3.

Here, our input is $x > 0$ and our solution is $f(x) = 3x$.

Our solution map is differentiable, $J(x) = f'(x) = 3$. Therefore the (relative) condition number is

$$\kappa = \frac{\|J(x)\|}{\|f(x)\|/\|x\|} = \frac{3}{|3x|/|x|} = 1.$$

Since κ is not large, the problem is well-conditioned (for all x).

Condition Number: Examples

Problem

Given a positive number, calculate its square root.

Here, our input is $x > 0$ and our solution is $f(x) = \sqrt{x}$.

Our solution map is differentiable, $J(x) = f'(x) = 1/(2\sqrt{x})$. Therefore,

$$\kappa = \frac{\|J(x)\|}{\|f(x)\|/\|x\|} = \frac{|1/(2\sqrt{x})|}{|\sqrt{x}|/|x|} = \frac{1}{2}.$$

Since κ is not large, the problem is well-conditioned (for all x).

Here, the absolute condition number $\tilde{\kappa} = \|J(x)\|$ can get very large for $x \approx 0$. However, we are always well-conditioned in terms of relative changes.

Condition Number: Examples

Problem

Given two numbers, calculate their difference.

Our input is $\mathbf{x} \in \mathbb{R}^2$ and our solution is $f(\mathbf{x}) = x_1 - x_2$. For convenience, let's use the ∞ -norm on \mathbb{R}^2 .

The Jacobian matrix is $J(\mathbf{x}) = \begin{bmatrix} \frac{\partial f}{\partial x_1} & \frac{\partial f}{\partial x_2} \end{bmatrix} = \begin{bmatrix} 1 & -1 \end{bmatrix}$. The operator ∞ -norm of this matrix is the maximum absolute row sum: $\|J(\mathbf{x})\|_\infty = |1| + |-1| = 2$.

Then our condition number is

$$\kappa = \frac{\|J(\mathbf{x})\|_\infty}{\|f(\mathbf{x})\|_\infty / \|\mathbf{x}\|_\infty} = \frac{2}{|x_1 - x_2| / \max(|x_1|, |x_2|)} = \frac{2 \max(|x_1|, |x_2|)}{|x_1 - x_2|}.$$

If $|x_1 - x_2|$ is small compared to $|x_1|$ or $|x_2|$, then the problem is ill-conditioned.

This is exactly when catastrophic cancellation occurs! (see floating point lectures)

Wilkinson's Polynomial

A famous early example of ill-conditioning is for finding the roots of polynomials.

Problem

Given $p(x) = a_0 + a_1x + \cdots + a_{n-1}x^{n-1} + x^n$, find the roots of $p(x)$.

Wilkinson's polynomial (James Wilkinson, 1963) is the simple-looking polynomial

$$p(x) = (x - 1)(x - 2) \cdots (x - 20).$$

Clearly $p(x)$ has 20 roots at $x = 1, 2, \dots, 20$.

Wilkinson wanted to test his implementation of Newton's method for rootfinding, so he ran it on the monomial form of $p(x)$:

$$p(x) = x^{20} - 210x^{19} + \cdots$$

Despite his code being correct, he could not find the roots!

Wilkinson's Polynomial

Suppose we want to find the 15th root ($x_{15} = 15$) as a function of the coefficient a_{15} .

With some work, it can be shown that the condition number is huge:

$$\kappa = \frac{a_{15} 15^{14}}{5! 14!} \approx 5.1 \times 10^{13}.$$

Suppose we are working in double precision, where $\epsilon \approx 2.22 \times 10^{-16}$. Assuming the relative error in the stored value of a_{15} is ϵ , then the relative error in the root is

$$\frac{|\delta x_{15}|}{|x_{15}|} \approx \kappa \frac{|\delta a_{15}|}{|a_{15}|} \leq \kappa \epsilon \approx 1.13 \times 10^{-2}.$$

That is, 16 correct digits in a_{15} leads to only 2 correct digits in x_{15} , regardless of how good our rootfinding algorithm is!

(Wilkinson was only working in single precision, where $\epsilon \approx 1.1 \times 10^{-7}$, so he got no digits of accuracy in his solution)

Condition Number: Examples

An important example of conditioning is matrix-vector multiplication.

Problem

Given $\mathbf{x} \in \mathbb{R}^n$ and $A \in \mathbb{R}^{n \times n}$, compute $A\mathbf{x}$.

Here, the solution is $\mathbf{y} = A\mathbf{x}$, and the perturbed solution is

$$\mathbf{y} + \delta\mathbf{y} = A(\mathbf{x} + \delta\mathbf{x}) = A\mathbf{x} + A\delta\mathbf{x},$$

and so $\delta\mathbf{y} = A\delta\mathbf{x}$.

The condition number for this problem is

$$\kappa = \max_{\delta\mathbf{x}} \frac{\|\delta\mathbf{y}\|/\|\mathbf{y}\|}{\|\delta\mathbf{x}\|/\|\mathbf{x}\|} = \max_{\delta\mathbf{x}} \frac{\|A\delta\mathbf{x}\|/\|A\mathbf{x}\|}{\|\delta\mathbf{x}\|/\|\mathbf{x}\|}$$

Condition Number: Examples

Noting that $\|A\delta\mathbf{x}\| \leq \|A\| \cdot \|\delta\mathbf{x}\|$, we get

$$\kappa \leq \frac{\|A\| \|\mathbf{x}\|}{\|A\mathbf{x}\|}.$$

If A is invertible, then $\frac{\|\mathbf{x}\|}{\|A\mathbf{x}\|} = \frac{\|A^{-1}(A\mathbf{x})\|}{\|A\mathbf{x}\|} \leq \|A^{-1}\|$, and so

$$\kappa \leq \|A\| \cdot \|A^{-1}\|.$$

Definition (Condition number of a matrix)

The condition number of an invertible matrix is

$$\kappa(A) = \|A\| \cdot \|A^{-1}\|.$$

This number appears in the analysis of many linear algebra problems.

Stability

Conditioning tells us how sensitive our mathematical problem is to perturbations in the input space.

In practice, we cannot evaluate the true solution $\mathbf{f}(\mathbf{x})$, but instead run some algorithm (implemented on a computer) which approximately finds the solution, $\tilde{\mathbf{f}}(\mathbf{x})$.

We hope that our algorithm is good, for example as measured by absolute or relative errors

$$\|\tilde{\mathbf{f}}(\mathbf{x}) - \mathbf{f}(\mathbf{x})\| \quad \text{or} \quad \frac{\|\tilde{\mathbf{f}}(\mathbf{x}) - \mathbf{f}(\mathbf{x})\|}{\|\mathbf{f}(\mathbf{x})\|}.$$

No matter the problem, there will always be some error in $\tilde{\mathbf{f}}(\mathbf{x})$, even if it is just from storing the data \mathbf{x} in finite precision.

If the problem is ill-conditioned (like Wilkinson's polynomial), this can be enough to cause large errors in $\tilde{\mathbf{f}}(\mathbf{x})$ regardless of the algorithm! So, we need some notion of a “good algorithm” that is achievable.

Definition (Stable algorithm)

An algorithm \tilde{f} is stable if, for all data $\mathbf{x} \in X$,

$$\frac{\|\tilde{f}(\mathbf{x}) - f(\tilde{\mathbf{x}})\|}{\|f(\tilde{\mathbf{x}})\|} = \mathcal{O}(\epsilon),$$

for some $\tilde{\mathbf{x}}$ with

$$\frac{\|\tilde{\mathbf{x}} - \mathbf{x}\|}{\|\mathbf{x}\|} = \mathcal{O}(\epsilon).$$

“A stable algorithm gives nearly the right answer to nearly the right question.” (Nick Trefethen & David Bau)

In practice, we can often do a bit better than this.

Definition (Backward stable algorithm)

An algorithm \tilde{f} is **backward stable** if, for all data $\mathbf{x} \in X$,

$$\tilde{f}(\mathbf{x}) = f(\tilde{\mathbf{x}}),$$

for some $\tilde{\mathbf{x}}$ with

$$\frac{\|\tilde{\mathbf{x}} - \mathbf{x}\|}{\|\mathbf{x}\|} = \mathcal{O}(\epsilon).$$

“A backwards stable algorithm gives **exactly** the right answer to nearly the right question.” (Nick Trefethen & David Bau)

Stability: Examples

All standard operations $+$, $-$, \times , \div are backwards stable. Let's consider subtraction.

Algorithm

Subtract two numbers using the algorithm

$$\tilde{f}(x_1, x_2) = \text{fl}(\text{fl}(x_1) - \text{fl}(x_2)).$$

$$\begin{aligned}\tilde{f}(x_1, x_2) &= \text{fl}(x_1(1 + \delta_1) - x_2(1 + \delta_2)), \\ &= [x_1(1 + \delta_1) - x_2(1 + \delta_2)](1 + \delta_3), \\ &= x_1 \underbrace{(1 + \delta_1)(1 + \delta_3)}_{=1+\delta_4} - x_2 \underbrace{(1 + \delta_2)(1 + \delta_3)}_{=1+\delta_5}.\end{aligned}$$

So $\tilde{f}(x_1, x_2) = \tilde{x}_1 - \tilde{x}_2 = f(\tilde{x}_1, \tilde{x}_2)$ where $\tilde{x}_1 = x_1(1 + \delta_4)$ and $\tilde{x}_2 = x_2(1 + \delta_5)$.

Since $|\delta_4|, |\delta_5| \leq 2\epsilon + \mathcal{O}(\epsilon^2)$, subtraction is backwards stable.

Algorithm

Calculate the eigenvalues of a matrix by

1. Calculate the characteristic polynomial and expand it into the form
$$p(x) = a_0 + a_1x + \cdots + a_nx^n.$$
2. Use a rootfinding algorithm to calculate the roots of $p(x)$.

The first step will introduce (small) rounding errors in the coefficients a_i of $p(x)$.

We know from Wilkinson's polynomial that small errors in the a_i can introduce large errors in the roots (since the problem is ill-conditioned).

Therefore the overall algorithm is not stable.

Stability: Examples

For example, suppose we want to find the eigenvalues of

$$A = \begin{bmatrix} 1 + 10^{-14} & 0 \\ 0 & 1 \end{bmatrix}.$$

The characteristic polynomial is

$$p(x) = \det(xI - A) = x^2 + (-2 - 10^{-14})x + (1 + 10^{-14}).$$

Using NumPy's polynomial rootfinding function `numpy.roots`, I get

$$x = 0.9999999850988439 \quad \text{and} \quad x = 1.0000000149011663$$

Despite the coefficients of $p(x)$ being exactly correct to machine precision (16 digits), I only get about 8 digits of accuracy in the eigenvalues (compare: `numpy.linalg.eig` gives 16 digits accuracy).

Why is backwards stability useful?

Well-conditioned problem + backwards stable algorithm = small relative error

To be more formal:

Theorem

Suppose a backwards stable algorithm $\tilde{\mathbf{f}}$ is used to solve a problem \mathbf{f} with input data \mathbf{x} and condition number $\kappa(\mathbf{x})$. Then for machine precision ϵ ,

$$\frac{\|\tilde{\mathbf{f}}(\mathbf{x}) - \mathbf{f}(\mathbf{x})\|}{\|\mathbf{f}(\mathbf{x})\|} = \mathcal{O}(\kappa(\mathbf{x})\epsilon).$$