

# An Introduction to Fixed-Point A Posteriori Validation

---

Florent Bréhard     [florent.brehard@univ-lille.fr](mailto:florent.brehard@univ-lille.fr)

Thursday, October 24, 2024

## Some Limitations of the Self-Validating Approach

Compute a **Rigorous Polynomial Approximation** of  $\frac{1}{P(x)}$

- Rational function  $\rightsquigarrow$  not a polynomial (of finite degree)

## Some Limitations of the Self-Validating Approach

Compute a **Rigorous Polynomial Approximation** of  $\frac{1}{P(x)}$

- Rational function  $\rightsquigarrow$  not a polynomial (of finite degree)
- **Taylor model:**  $P = (P_0 + P_1x + \cdots + P_nx^n, \varepsilon) \rightsquigarrow$   
 $Q = (Q_0 + Q_1x + \cdots + Q_nx^n, \eta)$   
 $P(x)Q(x) = 1 \Rightarrow P_n$  depends explicitly on  $P_0, \dots, P_n$  and  $Q_0, \dots, Q_{n-1}$   
How to determine  $\eta$  rigorously?

## Some Limitations of the Self-Validating Approach

Compute a **Rigorous Polynomial Approximation** of  $\frac{1}{P(x)}$

- Rational function  $\rightsquigarrow$  not a polynomial (of finite degree)
- **Taylor model:**  $P = (P_0 + P_1x + \cdots + P_nx^n, \varepsilon) \rightsquigarrow$   
 $Q = (Q_0 + Q_1x + \cdots + Q_nx^n, \eta)$   
 $P(x)Q(x) = 1 \Rightarrow P_n$  depends explicitly on  $P_0, \dots, P_n$  and  $Q_0, \dots, Q_{n-1}$   
How to determine  $\eta$  rigorously?
- **Chebyshev model:**  $P = (P_0 + P_1x + \cdots + P_nT_n(x), \varepsilon) \rightsquigarrow$   
 $Q = (Q_0 + Q_1x + \cdots + Q_nT_n(x), \eta)$   
 $P(x)Q(x) = 1 \Rightarrow$  not a finite formula, since  
 $T_n(x)T_m(x) = \frac{1}{2}(T_{n+m}(x) + T_{|n-m|}(x))$

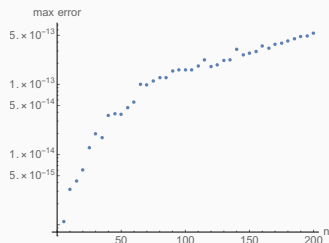
## Some Limitations of the Self-Validating Approach

Matrix inverse via Gaussian elimination using interval arithmetics

## Some Limitations of the Self-Validating Approach

Matrix inverse via Gaussian elimination using interval arithmetics

- **Example:** the Lehmer matrix  $L_n = \left( \frac{\min(i,j)}{\max(i,j)} \right)_{1 \leq i,j \leq n}$  is well-conditioned
- Gaussian elimination (using binary64 FP arithmetic) computes  $L_n^{-1}$  accurately

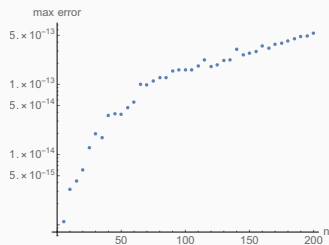


maximum error of  $L_n^{-1} L_n - I_n$  computed  
using binary64 Gaussian elimination

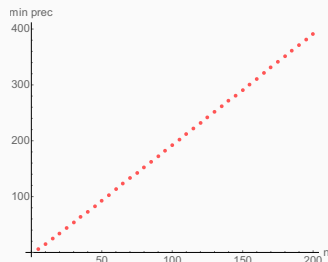
# Some Limitations of the Self-Validating Approach

Matrix inverse via Gaussian elimination using interval arithmetics

- **Example:** the Lehmer matrix  $L_n = \left( \frac{\min(i,j)}{\max(i,j)} \right)_{1 \leq i,j \leq n}$  is well-conditioned
- Gaussian elimination (using binary64 FP arithmetic) computes  $L_n^{-1}$  accurately
- Intervals of pivots using interval arithmetic grow much faster  
⇒ interval Gaussian elimination fails



maximum error of  $L_n^{-1} L_n - I_n$  computed using binary64 Gaussian elimination



minimum FP precision needed s.t. interval Gaussian elimination on  $L_n$  does not fail

## A Posteriori Validation Approach

- $\mathcal{F} : X \rightarrow Y$ ,  $X$  and  $Y$  Banach spaces, Solve  $\mathcal{F}(x) = 0 \rightsquigarrow x^* \in X$   
 $X = Y = (\mathcal{C}(I), \|\cdot\|_\infty)$ ,  $\mathcal{F}(f) = gf - 1 = 0$ ,  $f^* = \frac{1}{g}$



## A Posteriori Validation Approach

- $\mathcal{F} : X \rightarrow Y$ ,  $X$  and  $Y$  Banach spaces, Solve  $\mathcal{F}(x) = 0 \rightsquigarrow x^* \in X$   
 $X = Y = (\mathcal{C}(I), \|\cdot\|_\infty)$ ,  $\mathcal{F}(f) = gf - 1 = 0$ ,  $f^* = \frac{1}{g}$
- $\tilde{x} \in X$  a numerical solution:  $\mathcal{F}(\tilde{x}) \approx 0 \rightsquigarrow \tilde{x} \approx x^*$   
Compute  $\tilde{f} = p = \sum_{i=0}^n a_i T_i \in \mathbb{R}_n[x] \subseteq X$  by interpolating  $\frac{1}{g}$  at the Chebyshev nodes

## A Posteriori Validation Approach

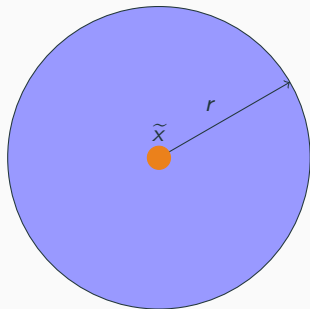
- $\mathcal{F} : X \rightarrow Y$ ,  $X$  and  $Y$  Banach spaces, Solve  $\mathcal{F}(x) = 0 \rightsquigarrow x^* \in X$   
 $X = Y = (\mathcal{C}(I), \|\cdot\|_\infty)$ ,  $\mathcal{F}(f) = gf - 1 = 0$ ,  $f^* = \frac{1}{g}$
- $\tilde{x} \in X$  a numerical solution:  $\mathcal{F}(\tilde{x}) \approx 0 \rightsquigarrow \tilde{x} \approx x^*$   
Compute  $\tilde{f} = p = \sum_{i=0}^n a_i T_i \in \mathbb{R}_n[x] \subseteq X$  by interpolating  $\frac{1}{g}$  at the Chebyshev nodes
- **A posteriori validation**  $\rightsquigarrow$  recover a rigorous bound  $\varepsilon \geq \|\tilde{x} - x^*\|$   
The pair  $(p, \varepsilon)$  is a Rigorous Polynomial Approximation for  $f^* = \frac{1}{g}$

## A Posteriori Validation Approach

- $\mathcal{F} : X \rightarrow Y$ ,  $X$  and  $Y$  Banach spaces, Solve  $\mathcal{F}(x) = 0 \rightsquigarrow x^* \in X$   
 $X = Y = (\mathcal{C}(I), \|\cdot\|_\infty)$ ,  $\mathcal{F}(f) = gf - 1 = 0$ ,  $f^* = \frac{1}{g}$
  - $\tilde{x} \in X$  a numerical solution:  $\mathcal{F}(\tilde{x}) \approx 0 \rightsquigarrow \tilde{x} \approx x^*$   
Compute  $\tilde{f} = p = \sum_{i=0}^n a_i T_i \in \mathbb{R}_n[x] \subseteq X$  by interpolating  $\frac{1}{g}$  at the Chebyshev nodes
  - **A posteriori validation**  $\rightsquigarrow$  recover a rigorous bound  $\varepsilon \geq \|\tilde{x} - x^*\|$   
The pair  $(p, \varepsilon)$  is a Rigorous Polynomial Approximation for  $f^* = \frac{1}{g}$
- $\Rightarrow$  Use a **fixed-point theorem!**

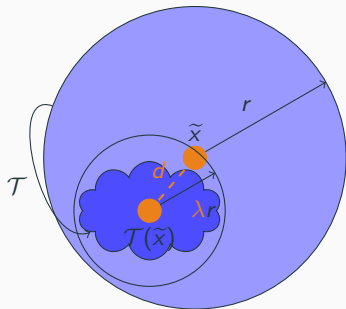
## Banach Fixed-Point Theorem

- Convert  $\mathcal{F}(x) = 0$  into an equivalent **fixed-point equation**  $\mathcal{T}(x) = x$   
 $\Rightarrow \mathcal{T} : X \rightarrow X$  must be **contracting**



# Banach Fixed-Point Theorem

- Convert  $\mathcal{F}(x) = 0$  into an equivalent **fixed-point equation**  $\mathcal{T}(x) = x$   
 $\Rightarrow \mathcal{T} : X \rightarrow X$  must be **contracting**



## Banach Fixed-Point Theorem

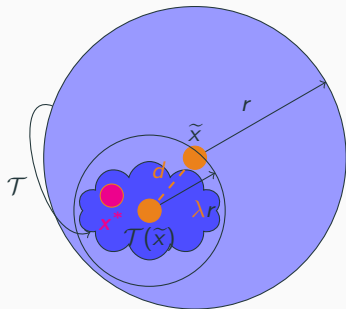
If one can rigorously check:

- $\mathcal{T}(B(\tilde{x}, r)) \subseteq B(\tilde{x}, r)$
- $\mathcal{T}$  is  $\lambda$ -contracting\* over  $B(\tilde{x}, r)$  with  $\lambda < 1$

\* it means that  $\|\mathcal{T}(x) - \mathcal{T}(x')\| \leq \lambda \|x - x'\|$  for all  $x, x' \in B(\tilde{x}, r)$

# Banach Fixed-Point Theorem

- Convert  $\mathcal{F}(x) = 0$  into an equivalent **fixed-point equation**  $\mathcal{T}(x) = x$   
 $\Rightarrow \mathcal{T} : X \rightarrow X$  must be **contracting**



## Banach Fixed-Point Theorem

If one can rigorously check:

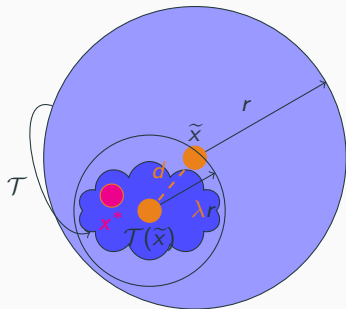
- $\mathcal{T}(B(\tilde{x}, r)) \subseteq B(\tilde{x}, r)$
- $\mathcal{T}$  is  $\lambda$ -contracting\* over  $B(\tilde{x}, r)$  with  $\lambda < 1$

Then  $\mathcal{T}$  has a unique fixed-point  $x^*$  in  $B(\tilde{x}, r)$

\* it means that  $\|\mathcal{T}(x) - \mathcal{T}(x')\| \leq \lambda \|x - x'\|$  for all  $x, x' \in B(\tilde{x}, r)$

# Banach Fixed-Point Theorem

- Convert  $\mathcal{F}(x) = 0$  into an equivalent **fixed-point equation**  $\mathcal{T}(x) = x$   
 $\Rightarrow \mathcal{T} : X \rightarrow X$  must be **contracting**



## Banach Fixed-Point Theorem

If one can rigorously check:

- $\mathcal{T}(B(\tilde{x}, r)) \subseteq B(\tilde{x}, r)$
- $\mathcal{T}$  is  $\lambda$ -contracting\* over  $B(\tilde{x}, r)$  with  $\lambda < 1$

Then  $\mathcal{T}$  has a unique fixed-point  $x^*$  in  $B(\tilde{x}, r)$

- Find “optimal” radius  $r$  that satisfies the theorem

\* it means that  $\|\mathcal{T}(x) - \mathcal{T}(x')\| \leq \lambda\|x - x'\|$  for all  $x, x' \in B(\tilde{x}, r)$

## A Linear Example: Division of RPAs (1/3)

Compute  $f = \frac{g}{h}$ ,  $h(x) \neq 0$  over  $I$ :

- Numerical interpolation  $\rightsquigarrow \tilde{f} \approx \frac{g}{h}$  in  $\mathbb{R}_n[x] \subseteq X = (\mathcal{C}(I), \|\cdot\|_\infty)$



## A Linear Example: Division of RPAs (1/3)

Compute  $f = \frac{g}{h}$ ,  $h(x) \neq 0$  over  $I$ :

- Numerical interpolation  $\rightsquigarrow \tilde{f} \approx \frac{g}{h}$  in  $\mathbb{R}_n[x] \subseteq X = (\mathcal{C}(I), \|\cdot\|_\infty)$
- We want to solve:

$$\mathcal{F}(f) = g \quad \text{where} \quad \mathcal{F} : X \rightarrow X, f \mapsto hf$$

$\Rightarrow \mathcal{F}$  is linear

## A Linear Example: Division of RPAs (1/3)

Compute  $f = \frac{g}{h}$ ,  $h(x) \neq 0$  over  $I$ :

- Numerical interpolation  $\rightsquigarrow \tilde{f} \approx \frac{g}{h}$  in  $\mathbb{R}_n[x] \subseteq X = (\mathcal{C}(I), \|\cdot\|_\infty)$
- We want to solve:

$$\mathcal{F}(f) = g \quad \text{where} \quad \mathcal{F} : X \rightarrow X, f \mapsto hf$$

$\Rightarrow \mathcal{F}$  is **linear**

- Since  $\mathcal{F}$  is linear, it coincides with its Fréchet derivative  $D\mathcal{F}_f : X \rightarrow X$

$$\mathcal{F}(f + \delta_f) = h(f + \delta_f) = \underbrace{hf}_{\mathcal{F}(f)} + \underbrace{h\delta_f}_{D\mathcal{F}_f(\delta_f)} \rightsquigarrow D\mathcal{F}_f(\delta_f) = \mathcal{F}(\delta_f) = h\delta_f$$

## A Linear Example: Division of RPAs (2/3)

Compute  $f = \frac{g}{h}$ ,  $h(x) \neq 0$  over  $I$ :

- Construct a **Newton-like fixed-point operator**  $\mathcal{T} : X \rightarrow X$ :

$$\mathcal{T}(f) = f - \mathcal{A}(\mathcal{F}(f) - g) \quad \text{where } \mathcal{A} \approx D\mathcal{F}_f^{-1} : X \rightarrow X$$

## A Linear Example: Division of RPAs (2/3)

Compute  $f = \frac{g}{h}$ ,  $h(x) \neq 0$  over  $I$ :

- Construct a **Newton-like fixed-point operator**  $\mathcal{T} : X \rightarrow X$ :

$$\mathcal{T}(f) = f - \mathcal{A}(\mathcal{F}(f) - g) \quad \text{where } \mathcal{A} \approx D\mathcal{F}_f^{-1} : X \rightarrow X$$

- $D\mathcal{F}_f^{-1} : \delta_f \mapsto \frac{\delta_f}{h}$ , so we define  $\mathcal{A}(\delta_f) = \tilde{\varphi}\delta_f$  using  $\tilde{\varphi} \approx \frac{1}{h}$

$$\mathcal{T}(f) = f - \tilde{\varphi}(hf - g) \quad \rightsquigarrow \quad \mathcal{T}(f) = f \Leftrightarrow f = \frac{g}{h}$$

## A Linear Example: Division of RPAs (2/3)

Compute  $f = \frac{g}{h}$ ,  $h(x) \neq 0$  over  $I$ :

- Construct a **Newton-like fixed-point operator**  $\mathcal{T} : X \rightarrow X$ :

$$\mathcal{T}(f) = f - \mathcal{A}(\mathcal{F}(f) - g) \quad \text{where} \quad \mathcal{A} \approx D\mathcal{F}_f^{-1} : X \rightarrow X$$

- $D\mathcal{F}_f^{-1} : \delta_f \mapsto \frac{\delta_f}{h}$ , so we define  $\mathcal{A}(\delta_f) = \tilde{\varphi}\delta_f$  using  $\tilde{\varphi} \approx \frac{1}{h}$

$$\mathcal{T}(f) = f - \tilde{\varphi}(hf - g) \quad \rightsquigarrow \quad \mathcal{T}(f) = f \Leftrightarrow f = \frac{g}{h}$$

- Check the contraction property of  $\mathcal{T}$ :

$$\begin{aligned} \|\mathcal{T}(f) - \mathcal{T}(f')\| &= \|[f - \tilde{\varphi}(hf - g)] - [f' - \tilde{\varphi}(hf' - g)]\| \\ &= \|(1 - \tilde{\varphi}h)(f - f')\| \leq \underbrace{\|1 - \tilde{\varphi}h\|}_{:=\lambda} \|f - f'\| \end{aligned}$$

## A Linear Example: Division of RPAs (2/3)

Compute  $f = \frac{g}{h}$ ,  $h(x) \neq 0$  over  $I$ :

- Construct a **Newton-like fixed-point operator**  $\mathcal{T} : X \rightarrow X$ :

$$\mathcal{T}(f) = f - \mathcal{A}(\mathcal{F}(f) - g) \quad \text{where} \quad \mathcal{A} \approx D\mathcal{F}_f^{-1} : X \rightarrow X$$

- $D\mathcal{F}_f^{-1} : \delta_f \mapsto \frac{\delta_f}{h}$ , so we define  $\mathcal{A}(\delta_f) = \tilde{\varphi}\delta_f$  using  $\tilde{\varphi} \approx \frac{1}{h}$

$$\mathcal{T}(f) = f - \tilde{\varphi}(hf - g) \quad \rightsquigarrow \quad \mathcal{T}(f) = f \Leftrightarrow f = \frac{g}{h}$$

- Check the contraction property of  $\mathcal{T}$ :

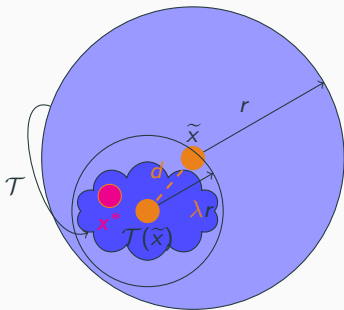
$$\begin{aligned} \|\mathcal{T}(f) - \mathcal{T}(f')\| &= \|[f - \tilde{\varphi}(hf - g)] - [f' - \tilde{\varphi}(hf' - g)]\| \\ &= \|(1 - \tilde{\varphi}h)(f - f')\| \leq \underbrace{\|1 - \tilde{\varphi}h\|}_{:=\lambda} \|f - f'\| \end{aligned}$$

$\Rightarrow$  We need  $\lambda = \|1 - \tilde{\varphi}h\| < 1$

## A Linear Example: Division of RPAs (3/3)

Compute  $f = \frac{g}{h}$ ,  $h(x) \neq 0$  over  $I$ :

- Check the stability condition to apply the Banach fixed-point theorem



### Banach Fixed-Point Theorem

If one can rigorously check:

- $\mathcal{T}(B(\tilde{f}, r)) \subseteq B(\tilde{f}, r)$
- $\mathcal{T}$  is  $\lambda$ -contracting\* over  $B(\tilde{x}, r)$  with  $\lambda < 1$

Then  $\mathcal{T}$  has a unique fixed-point  $x^*$  in  $B(\tilde{x}, r)$

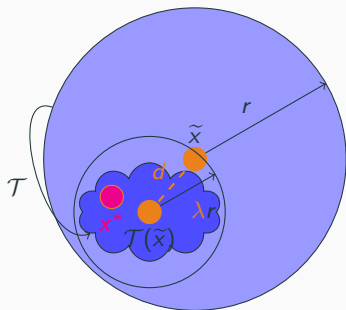
- The stability condition is encoded as:

$$d + \lambda r \leq r, \quad d = \|\mathcal{T}(\tilde{f}) - \tilde{f}\| = \|\tilde{\varphi}(h\tilde{f} - g)\|$$

## A Linear Example: Division of RPAs (3/3)

Compute  $f = \frac{g}{h}$ ,  $h(x) \neq 0$  over  $I$ :

- Check the stability condition to apply the Banach fixed-point theorem



### Banach Fixed-Point Theorem

If one can rigorously check:

- $\mathcal{T}(B(\tilde{f}, r)) \subseteq B(\tilde{f}, r)$
- $\mathcal{T}$  is  $\lambda$ -contracting\* over  $B(\tilde{x}, r)$  with  $\lambda < 1$

Then  $\mathcal{T}$  has a unique fixed-point  $x^*$  in  $B(\tilde{x}, r)$

- The stability condition is encoded as:

$$d + \lambda r \leq r, \quad d = \|\mathcal{T}(\tilde{f}) - \tilde{f}\| = \|\tilde{\varphi}(h\tilde{f} - g)\|$$

$\Rightarrow$  the “best” bound  $r$  is  $r = \frac{d}{1 - \lambda}$



### Algorithm RPADiv

- **Input:** RPAs  $g$  and  $h$ , approximation degree  $n \in \mathbb{N}$
- **Output:** degree- $n$  RPA  $f$  representing  $\frac{g}{h}$  rigorously

1. Compute  $\tilde{f} \approx \frac{g}{h}$  using degree- $n$  Chebyshev interpolation
2. Compute  $\tilde{\varphi} \approx \frac{1}{h}$  using degree- $n$  Chebyshev interpolation
3. Compute  $\lambda = \|1 - \tilde{\varphi}h\|$  and **FAIL** if  $\lambda \geq 1$
4. Compute  $d = \|\tilde{\varphi}(h\tilde{f} - g)\|$
5. Compute  $r = \frac{d}{1 - \lambda}$  and **RETURN** RPA  $f = (\tilde{f}, r)$

### Algorithm RPADiv is correct

If  $\text{RPADiv}(g, h)$  does not fail, then it returns an RPA  $f$  such that for all  $g \in \mathbf{g}, h \in \mathbf{h}$ , we have  $f = \frac{g}{h} \in \mathbf{f}$ .

## A Nonlinear Example: Square Root of an RPA (1/3)

Compute  $f = \sqrt{g}$ ,  $g(x) > 0$  over  $I$ :

- Numerical interpolation  $\rightsquigarrow \tilde{f} \approx \sqrt{g}$  in  $\mathbb{R}_n[x] \subseteq X = (\mathcal{C}(I), \|\cdot\|_\infty)$

## A Nonlinear Example: Square Root of an RPA (1/3)

Compute  $f = \sqrt{g}$ ,  $g(x) > 0$  over  $I$ :

- Numerical interpolation  $\rightsquigarrow \tilde{f} \approx \sqrt{g}$  in  $\mathbb{R}_n[x] \subseteq X = (C(I), \|\cdot\|_\infty)$
- We want to solve:

$$\mathcal{F}(f) = g \quad \text{where} \quad \mathcal{F} : X \rightarrow X, f \mapsto f^2$$

$\Rightarrow \mathcal{F}$  is **nonlinear** (quadratic)

## A Nonlinear Example: Square Root of an RPA (1/3)

Compute  $f = \sqrt{g}$ ,  $g(x) > 0$  over  $I$ :

- Numerical interpolation  $\rightsquigarrow \tilde{f} \approx \sqrt{g}$  in  $\mathbb{R}_n[x] \subseteq X = (\mathcal{C}(I), \|\cdot\|_\infty)$
- We want to solve:

$$\mathcal{F}(f) = g \quad \text{where} \quad \mathcal{F} : X \rightarrow X, f \mapsto f^2$$

$\Rightarrow \mathcal{F}$  is **nonlinear** (quadratic)

- Since  $\mathcal{F}$  is nonlinear, its Fréchet derivative  $D\mathcal{F}_f : X \rightarrow X$  **depends** on  $f$ :

$$\mathcal{F}(f + \delta_f) = (f + \delta_f)^2 = \underbrace{f^2}_{\mathcal{F}(f)} + \underbrace{2f\delta_f}_{D\mathcal{F}_f(\delta_f)} + \delta_f^2 \quad \rightsquigarrow \quad D\mathcal{F}_f(\delta_f) = 2f\delta_f$$

## A Nonlinear Example: Square Root of an RPA (2/3)

Compute  $f = \sqrt{g}$ ,  $g(x) > 0$  over  $I$ :

- Construct a **Newton-like fixed-point operator**  $\mathcal{T} : X \rightarrow X$ :

$$\mathcal{T}(f) = f - \mathcal{A}(\mathcal{F}(f) - g) \quad \text{where } \mathcal{A} \approx D\mathcal{F}_f^{-1} : X \rightarrow X$$

## A Nonlinear Example: Square Root of an RPA (2/3)

Compute  $f = \sqrt{g}$ ,  $g(x) > 0$  over  $I$ :

- Construct a **Newton-like fixed-point operator**  $\mathcal{T} : X \rightarrow X$ :

$$\mathcal{T}(f) = f - \mathcal{A}(\mathcal{F}(f) - g) \quad \text{where } \mathcal{A} \approx D\mathcal{F}_f^{-1} : X \rightarrow X$$

- $D\mathcal{F}_f^{-1} : \delta_f \mapsto \frac{\delta_f}{2f}$ , so we define  $\mathcal{A}(\delta_f) = \tilde{\varphi}\delta_f$  using  $\tilde{\varphi} \approx \frac{1}{2f} \approx \frac{1}{2\sqrt{g}}$

$$\mathcal{T}(f) = f - \tilde{\varphi}(f^2 - g) \quad \rightsquigarrow \quad \mathcal{T}(f) = f \Leftrightarrow f = \pm\sqrt{g}$$

## A Nonlinear Example: Square Root of an RPA (2/3)

Compute  $f = \sqrt{g}$ ,  $g(x) > 0$  over  $I$ :

- Construct a **Newton-like fixed-point operator**  $\mathcal{T} : X \rightarrow X$ :

$$\mathcal{T}(f) = f - \mathcal{A}(\mathcal{F}(f) - g) \quad \text{where} \quad \mathcal{A} \approx D\mathcal{F}_f^{-1} : X \rightarrow X$$

- $D\mathcal{F}_f^{-1} : \delta_f \mapsto \frac{\delta_f}{2f}$ , so we define  $\mathcal{A}(\delta_f) = \tilde{\varphi}\delta_f$  using  $\tilde{\varphi} \approx \frac{1}{2f} \approx \frac{1}{2\sqrt{g}}$

$$\mathcal{T}(f) = f - \tilde{\varphi}(f^2 - g) \quad \rightsquigarrow \quad \mathcal{T}(f) = f \Leftrightarrow f = \pm\sqrt{g}$$

- Check the contraction property of  $\mathcal{T}$ , for  $f, f' \in B(\tilde{r})$ :

$$\begin{aligned} \|\mathcal{T}(f) - \mathcal{T}(f')\| &= \|[f - \tilde{\varphi}(f^2 - g)] - [f' - \tilde{\varphi}(f'^2 - g)]\| \\ &= \|[1 - \tilde{\varphi}(f + f')](f - f')\| \leq \underbrace{\|1 - \tilde{\varphi}(f + f')\|}_{\leq \lambda(r)} \|f - f'\| \end{aligned}$$

## A Nonlinear Example: Square Root of an RPA (2/3)

Compute  $f = \sqrt{g}$ ,  $g(x) > 0$  over  $I$ :

- Construct a **Newton-like fixed-point operator**  $\mathcal{T} : X \rightarrow X$ :

$$\mathcal{T}(f) = f - \mathcal{A}(\mathcal{F}(f) - g) \quad \text{where } \mathcal{A} \approx D\mathcal{F}_f^{-1} : X \rightarrow X$$

- $D\mathcal{F}_f^{-1} : \delta_f \mapsto \frac{\delta_f}{2f}$ , so we define  $\mathcal{A}(\delta_f) = \tilde{\varphi}\delta_f$  using  $\tilde{\varphi} \approx \frac{1}{2f} \approx \frac{1}{2\sqrt{g}}$

$$\mathcal{T}(f) = f - \tilde{\varphi}(f^2 - g) \quad \rightsquigarrow \quad \mathcal{T}(f) = f \Leftrightarrow f = \pm\sqrt{g}$$

- Check the contraction property of  $\mathcal{T}$ , for  $f, f' \in B(\tilde{r})$ :

$$\begin{aligned} \|\mathcal{T}(f) - \mathcal{T}(f')\| &= \|[f - \tilde{\varphi}(f^2 - g)] - [f' - \tilde{\varphi}(f'^2 - g)]\| \\ &= \|[1 - \tilde{\varphi}(f + f')](f - f')\| \leq \underbrace{\|1 - \tilde{\varphi}(f + f')\|}_{\leq \lambda(r)} \|f - f'\| \end{aligned}$$

$\Rightarrow \mathcal{T}$  is  $\lambda(r)$  **contracting** over  $B(\tilde{f}, r)$ , where:

$$\lambda(r) := \underbrace{\|1 - 2\tilde{\varphi}\tilde{f}\|}_{\lambda_0} + \underbrace{2\|\tilde{\varphi}\|}_{\lambda_1} r$$



## A Nonlinear Example: Square Root of an RPA (2/3)

Compute  $f = \sqrt{g}$ ,  $g(x) > 0$  over  $I$ :

- Construct a **Newton-like fixed-point operator**  $\mathcal{T} : X \rightarrow X$ :

$$\mathcal{T}(f) = f - \mathcal{A}(\mathcal{F}(f) - g) \quad \text{where } \mathcal{A} \approx D\mathcal{F}_f^{-1} : X \rightarrow X$$

- $D\mathcal{F}_f^{-1} : \delta_f \mapsto \frac{\delta_f}{2f}$ , so we define  $\mathcal{A}(\delta_f) = \tilde{\varphi}\delta_f$  using  $\tilde{\varphi} \approx \frac{1}{2f} \approx \frac{1}{2\sqrt{g}}$

$$\mathcal{T}(f) = f - \tilde{\varphi}(f^2 - g) \quad \rightsquigarrow \quad \mathcal{T}(f) = f \Leftrightarrow f = \pm\sqrt{g}$$

- Check the contraction property of  $\mathcal{T}$ , for  $f, f' \in B(\tilde{r})$ :

$$\begin{aligned} \|\mathcal{T}(f) - \mathcal{T}(f')\| &= \|[f - \tilde{\varphi}(f^2 - g)] - [f' - \tilde{\varphi}(f'^2 - g)]\| \\ &= \|[1 - \tilde{\varphi}(f + f')](f - f')\| \leq \underbrace{\|1 - \tilde{\varphi}(f + f')\|}_{\leq \lambda(r)} \|f - f'\| \end{aligned}$$

$\Rightarrow \mathcal{T}$  is  $\lambda(r)$  **contracting** over  $B(\tilde{r}, r)$ , where:

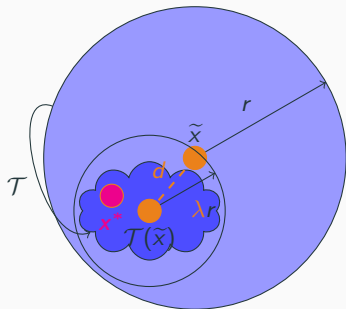
$$\lambda(r) := \underbrace{\|1 - 2\tilde{\varphi}\tilde{f}\|}_{\lambda_0} + \underbrace{2\|\tilde{\varphi}\|}_{\lambda_1} r$$

$\Rightarrow r$  must be small enough to ensure  $\lambda(r) \leq 1$

## A Nonlinear Example: Square Root of an RPA (3/3)

Compute  $f = \sqrt{g}$ ,  $g(x) > 0$  over  $I$ :

- Check the stability condition to apply the Banach fixed-point theorem



### Banach Fixed-Point Theorem

If one can rigorously check:

- $\mathcal{T}(B(\tilde{f}, r)) \subseteq B(\tilde{f}, r)$
- $\mathcal{T}$  is  $\lambda$ -contracting\* over  $B(\tilde{x}, r)$  with  $\lambda < 1$

Then  $\mathcal{T}$  has a unique fixed-point  $x^*$  in  $B(\tilde{x}, r)$

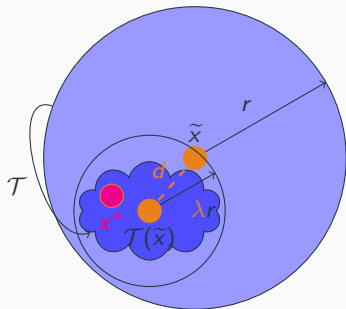
The stability condition is encoded as:

$$d + \lambda(r)r \leq r, \quad d = \|\mathcal{T}(\tilde{f}) - \tilde{f}\| = \|\tilde{\varphi}(\tilde{f}^2 - g)\|$$

## A Nonlinear Example: Square Root of an RPA (3/3)

Compute  $f = \sqrt{g}$ ,  $g(x) > 0$  over  $I$ :

- Check the stability condition to apply the Banach fixed-point theorem



### Banach Fixed-Point Theorem

If one can rigorously check:

- $\mathcal{T}(B(\tilde{f}, r)) \subseteq B(\tilde{f}, r)$
- $\mathcal{T}$  is  $\lambda$ -contracting\* over  $B(\tilde{x}, r)$  with  $\lambda < 1$

Then  $\mathcal{T}$  has a unique fixed-point  $x^*$  in  $B(\tilde{x}, r)$

The stability condition is encoded as:

$$d + \lambda_1 r \leq r, \quad d = \|\mathcal{T}(\tilde{f}) - \tilde{f}\| = \|\tilde{\varphi}(\tilde{f}^2 - g)\|$$

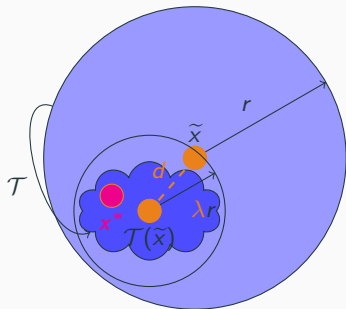
- $\lambda_1 r^2 + (\lambda_0 - 1)r + d \leq 0$  has positive real solutions iff:

$$\Delta := (1 - \lambda_0)^2 - 4\lambda_1 d \geq 0$$

## A Nonlinear Example: Square Root of an RPA (3/3)

Compute  $f = \sqrt{g}$ ,  $g(x) > 0$  over  $I$ :

- Check the stability condition to apply the Banach fixed-point theorem



### Banach Fixed-Point Theorem

If one can rigorously check:

- $\mathcal{T}(B(\tilde{f}, r)) \subseteq B(\tilde{f}, r)$
- $\mathcal{T}$  is  $\lambda$ -contracting\* over  $B(\tilde{x}, r)$  with  $\lambda < 1$

Then  $\mathcal{T}$  has a unique fixed-point  $x^*$  in  $B(\tilde{x}, r)$

The stability condition is encoded as:

$$d + \lambda(r)r \leq r, \quad d = \|\mathcal{T}(\tilde{f}) - \tilde{f}\| = \|\tilde{\varphi}(\tilde{f}^2 - g)\|$$

- $\lambda_1 r^2 + (\lambda_0 - 1)r + d \leq 0$  has positive real solutions iff:

$$\Delta := (1 - \lambda_0)^2 - 4\lambda_1 d \geq 0$$

$\Rightarrow$  Return the smallest root:  $r := \frac{1 - \lambda_0 - \sqrt{\Delta}}{2\lambda_1}$

## Square Root of an RPA — The Algorithm

### Algorithm RPASqrt

- **Input:** RPA  $\mathbf{g}$ , approximation degree  $n \in \mathbb{N}$
  - **Output:** degree- $n$  RPA  $\mathbf{f}$  representing  $\sqrt{\mathbf{g}}$  rigorously
- 
1. Compute  $\tilde{\mathbf{f}} \approx \sqrt{\mathbf{g}}$  using degree- $n$  Chebyshev interpolation
  2. Compute  $\tilde{\varphi} \approx \frac{1}{2\tilde{\mathbf{f}}}$  using degree- $n$  Chebyshev interpolation
  3. Compute  $\lambda_0 = \|1 - 2\tilde{\varphi}\tilde{\mathbf{f}}\|$  and **FAIL** if  $\lambda_0 \geq 1$
  4. Compute  $\lambda_1 = 2\|\tilde{\varphi}\|$
  5. Compute  $d = \|\tilde{\varphi}(\tilde{\mathbf{f}}^2 - \mathbf{g})\|$
  6. Compute  $\Delta = (1 - \lambda_0)^2 - 4\lambda_1 d$  and **FAIL** if  $\Delta < 0$
  5. Compute  $r := \frac{1 - \lambda_0 - \sqrt{\Delta}}{2\lambda_1}$  and **RETURN** RPA  $\mathbf{f} = (\tilde{\mathbf{f}}, r)$

### Algorithm RPADiv is correct

If  $\text{RPASqrt}(\mathbf{g})$  does not fail, then it returns an RPA  $\mathbf{f}$  such that for all  $g \in \mathbf{g}$ , we have  $f = \sqrt{g} \in \mathbf{f}$ .

## Further Examples of Fixed-Point A Posteriori Validation (1/2)

- Roots of univariate polynomials:

$$p(x) = 0 \rightsquigarrow z_1, \dots, z_n \text{ s.t. } p(x) = (x - z_1) \dots (x - z_n)$$

- Roots of univariate analytic functions:

$$f(x) = 0 \rightsquigarrow \text{isolate one/all root(s) of } f$$

- Roots of systems of multivariate polynomial/analytic functions:

$$\begin{cases} f_1(x_1, \dots, x_n) = 0 \\ \vdots \\ f_n(x_1, \dots, x_n) = 0 \end{cases} \rightsquigarrow \text{isolate one/all solutions}$$

See for instance: Siegfried M. Rump, Verification methods: Rigorous results using floating-point arithmetic. *Acta Numerica*. 2010;19:287-449.

## Further Examples of Fixed-Point A Posteriori Validation (2/2)

- Solving linear systems:

$$Ax = b \rightsquigarrow \text{recover } x = (x_1, \dots, x_n)$$

- Eigenvalue problems:

$$Av = \lambda v \rightsquigarrow \text{recover eigenvalue/eigenvector pair } (\lambda, v)$$

$$A \rightsquigarrow \text{diagonalize } A = PDP^{-1} : \begin{cases} D = \text{diag}(\lambda_1, \dots, \lambda_n) = \text{eigenvalues} \\ P \in \mathbb{C}^{n \times n} = \text{eigenvectors} \end{cases}$$

- Ordinary differential equations (ODEs):

$$\begin{cases} y'(x) = f(x, y(x)) \\ y(0) = v \in \mathbb{R}^N \end{cases} \rightsquigarrow \text{compute } \mathbf{p} = (p, \varepsilon) \text{ for } y$$

- Partial Differential Equations (PDEs), Delay Differential Equations (DDEs), etc.

See for instance: Siegfried M. Rump, Verification methods: Rigorous results using floating-point arithmetic. *Acta Numerica*. 2010;19:287-449.

# Rigorous Numerics for Hilbert's 16th Problem

---

Florent Bréhard, Nicolas Brisebarre, Mioara Joldes, Damien Pous, Warwick Tucker

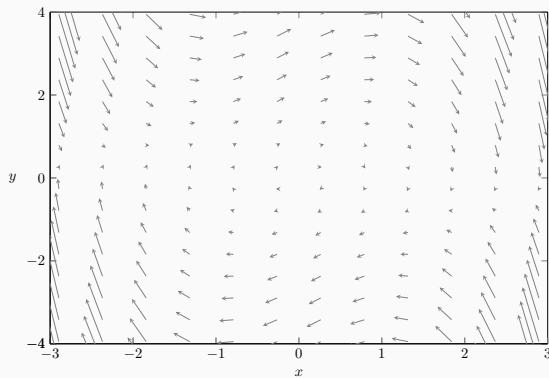
✉ [florent.brehard@univ-lille.fr](mailto:florent.brehard@univ-lille.fr)

Thursday, October 24, 2024



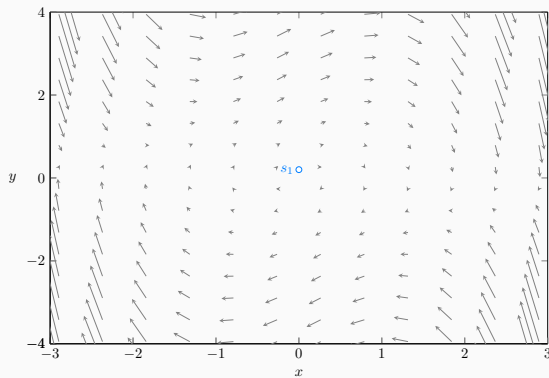
# Polynomial Vector Fields in the Plane

$$\begin{cases} \dot{x} = P(x, y) \\ \dot{y} = Q(x, y) \end{cases} \quad P, Q \in \mathbb{R}[x, y]$$



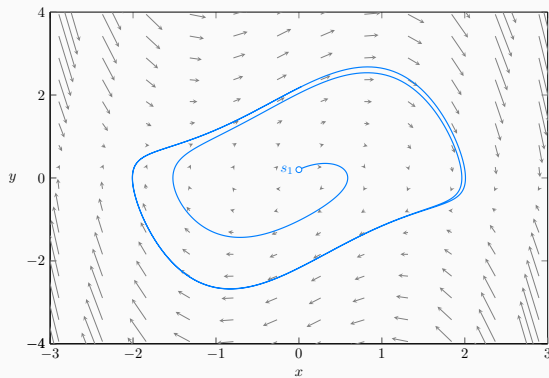
# Polynomial Vector Fields in the Plane

$$\begin{cases} \dot{x} = P(x, y) \\ \dot{y} = Q(x, y) \end{cases} \quad P, Q \in \mathbb{R}[x, y]$$



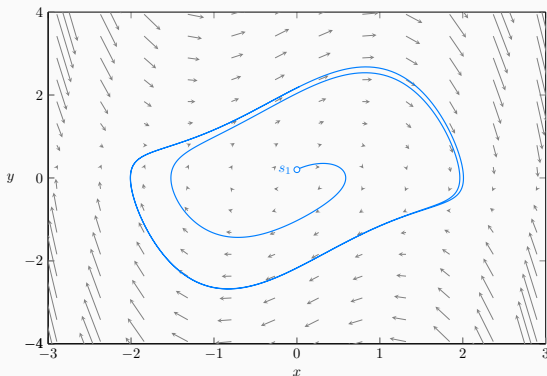
# Polynomial Vector Fields in the Plane

$$\begin{cases} \dot{x} = P(x, y) \\ \dot{y} = Q(x, y) \end{cases} \quad P, Q \in \mathbb{R}[x, y]$$



# Polynomial Vector Fields in the Plane

$$\begin{cases} \dot{x} = P(x, y) \\ \dot{y} = Q(x, y) \end{cases} \quad P, Q \in \mathbb{R}[x, y]$$



Examples:

- $(v, a) = (\dot{x}, \dot{v}) = (v, Q(x, v))$  in mechanics
- $(\dot{u}, \dot{i}) = (P(u, i), Q(u, i))$  in electricity

## Limit Cycles

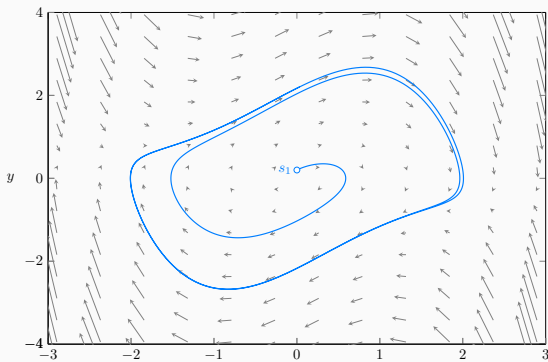
**limit cycle** = isolated periodic orbit

(For example,  $(\dot{x}, \dot{y}) = (-y, x)$  produces a continuum of periodic orbits but no limit cycles!)

# Limit Cycles

limit cycle = isolated periodic orbit

(For example,  $(\dot{x}, \dot{y}) = (-y, x)$  produces a continuum of periodic orbits but no limit cycles!)



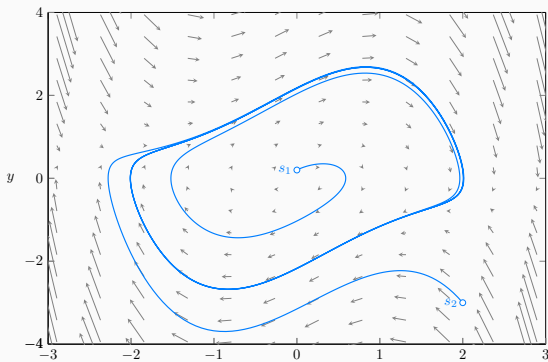
Van der Pol oscillator

$$\begin{cases} \dot{x} = y \\ \dot{y} = \mu(1 - x^2)y - x \end{cases}$$

# Limit Cycles

limit cycle = isolated periodic orbit

(For example,  $(\dot{x}, \dot{y}) = (-y, x)$  produces a continuum of periodic orbits but no limit cycles!)



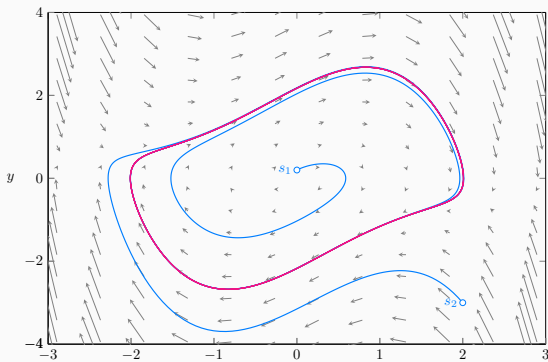
Van der Pol oscillator

$$\begin{cases} \dot{x} = y \\ \dot{y} = \mu(1 - x^2)y - x \end{cases}$$

# Limit Cycles

limit cycle = isolated periodic orbit

(For example,  $(\dot{x}, \dot{y}) = (-y, x)$  produces a continuum of periodic orbits but no limit cycles!)



Van der Pol oscillator 
$$\begin{cases} \dot{x} = y \\ \dot{y} = \mu(1 - x^2)y - x \end{cases}$$

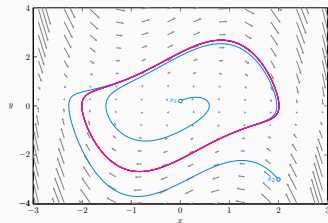


# Hilbert's 16th Problem

## Hilbert's 16th problem (second part)

For a given integer  $n$ , what is the maximum number  $\mathcal{H}(n)$  of **limit cycles** a **polynomial** vector field of degree **at most  $n$**  in the **plane** can have?

D. Hilbert, International Congress of Mathematicians, Paris, 1900



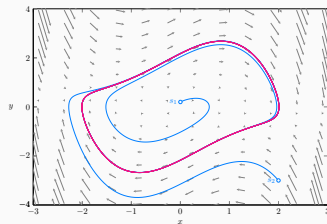
# Hilbert's 16th Problem

## Hilbert's 16th problem (second part)

For a given integer  $n$ , what is the maximum number  $\mathcal{H}(n)$  of **limit cycles** a **polynomial** vector field of degree **at most  $n$**  in the **plane** can have?

D. Hilbert, International Congress of Mathematicians, Paris, 1900

- 1923: H. Dulac (incorrectly) proved that a *single* polynomial vector field has a finite number of limit cycles
- 1981: Y. S. Il'Yashenko found a major gap in Dulac's proof
- 1991: New proofs of Dulac's result by Y. S. Il'Yashenko and J. Écalle
- But even  $\mathcal{H}(2) < \infty$  is open!

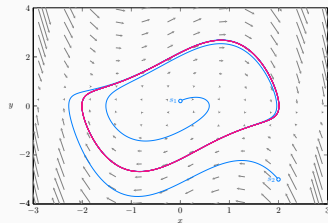


## Hilbert's 16th problem (second part)

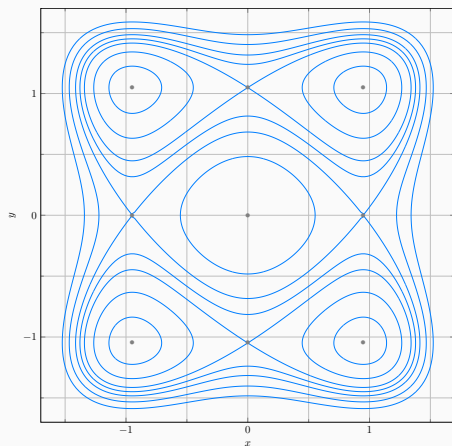
For a given integer  $n$ , what is the maximum number  $\mathcal{H}(n)$  of **limit cycles** a **polynomial** vector field of degree **at most  $n$**  in the **plane** can have?

D. Hilbert, International Congress of Mathematicians, Paris, 1900

- 1923: H. Dulac (incorrectly) proved that a *single* polynomial vector field has a finite number of limit cycles
- 1981: Y. S. Il'Yashenko found a major gap in Dulac's proof
- 1991: New proofs of Dulac's result by Y. S. Il'Yashenko and J. Écalle
- But even  $\mathcal{H}(2) < \infty$  is open!
- Some lower bounds:  $\mathcal{H}(2) \geq 4$ ,  $\mathcal{H}(3) \geq 13$ ,  $\mathcal{H}(4) \geq 28$ .



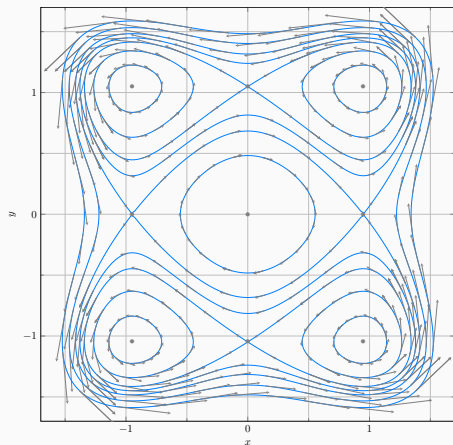
## Infinitesimal Hilbert's 16th Problem



$$H(x, y) = (x^2 - 0.9)^2 + (y^2 - 1.1)^2$$

T. Johnson, A quartic system with twenty-six limit cycles, *Experimental Mathematics*, 2011

# Infinitesimal Hilbert's 16th Problem

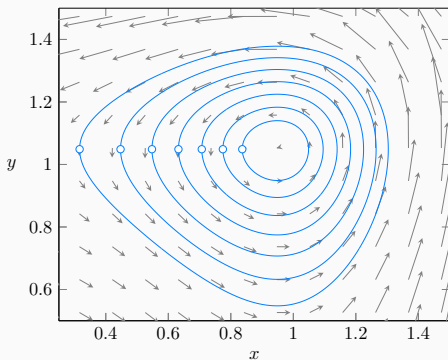


$$H(x, y) = (x^2 - 0.9)^2 + (y^2 - 1.1)^2$$

$$\begin{cases} \dot{x} = -\partial_y H(x, y) = 4y(y^2 - 1.1) \\ \dot{y} = \partial_x H(x, y) = 4x(x^2 - 0.9) \end{cases}$$

T. Johnson, A quartic system with twenty-six limit cycles, *Experimental Mathematics*, 2011

# Infinitesimal Hilbert's 16th Problem

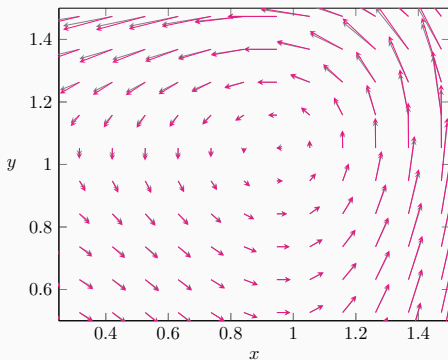


$$H(x, y) = (x^2 - 0.9)^2 + (y^2 - 1.1)^2$$

$$\begin{cases} \dot{x} = 4y(y^2 - 1.1) \\ \dot{y} = 4x(x^2 - 0.9) \end{cases}$$

T. Johnson, A quartic system with twenty-six limit cycles, *Experimental Mathematics*, 2011

# Infinitesimal Hilbert's 16th Problem

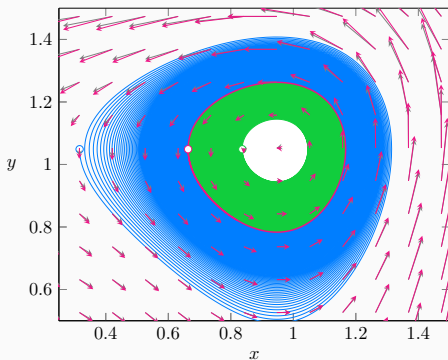


$$H(x, y) = (x^2 - 0.9)^2 + (y^2 - 1.1)^2$$

$$\begin{cases} \dot{x} = 4y(y^2 - 1.1) \\ \dot{y} = 4x(x^2 - 0.9) - 0.4y + 0.46x^2y \end{cases}$$

T. Johnson, A quartic system with twenty-six limit cycles, *Experimental Mathematics*, 2011

# Infinitesimal Hilbert's 16th Problem



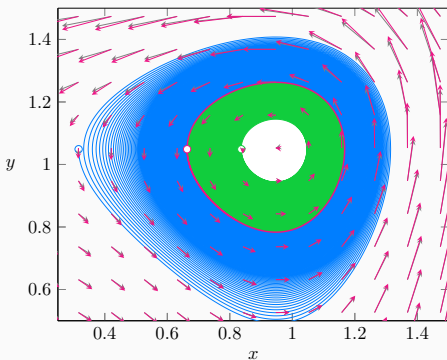
T. Johnson, A quartic system with twenty-six limit cycles, *Experimental Mathematics*, 2011

$$H(x, y) = (x^2 - 0.9)^2 + (y^2 - 1.1)^2$$

$$\begin{cases} \dot{x} = 4y(y^2 - 1.1) \\ \dot{y} = 4x(x^2 - 0.9) - 0.4y + 0.46x^2y \end{cases}$$



# Infinitesimal Hilbert's 16th Problem



T. Johnson, A quartic system with twenty-six limit cycles, *Experimental Mathematics*, 2011

## Infinitesimal Hilbert's 16th problem

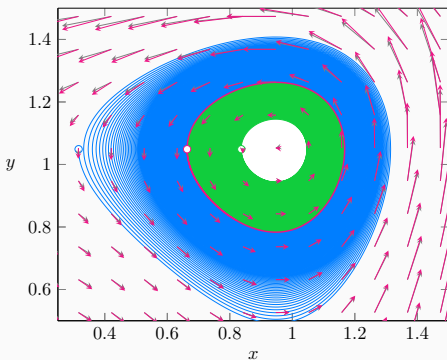
For a given integer  $n$ , what is the maximal number  $\mathcal{Z}(n)$  of limit cycles a **perturbed Hamiltonian** vector field of the form:

$$\begin{cases} \dot{x} = -\partial_y H(x, y) + \varepsilon f(x, y) \\ \dot{y} = \partial_x H(x, y) + \varepsilon g(x, y) \end{cases}$$

can have when  $\varepsilon \rightarrow 0$ , with:

- $H(x, y)$  a polynomial potential function of degree  $n + 1$
- $f, g$  polynomial perturbations of degree  $n$

# Infinitesimal Hilbert's 16th Problem



T. Johnson, A quartic system with twenty-six limit cycles, *Experimental Mathematics*, 2011

## Infinitesimal Hilbert's 16th problem

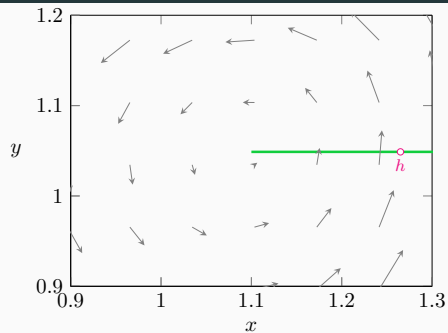
For a given integer  $n$ , what is the maximal number  $\mathcal{Z}(n)$  of limit cycles a **perturbed Hamiltonian** vector field of the form:

$$\begin{cases} \dot{x} = -\partial_y H(x, y) + \varepsilon f(x, y) \\ \dot{y} = \partial_x H(x, y) + \varepsilon g(x, y) \end{cases}$$

can have when  $\varepsilon \rightarrow 0$ , with:

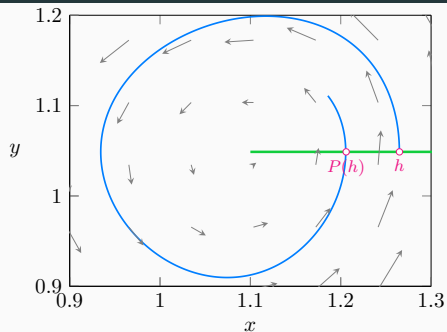
- $H(x, y)$  a polynomial potential function of degree  $n + 1$
  - $f, g$  polynomial perturbations of degree  $n$
  - $\mathcal{Z}(n) < \infty$  for all  $n$
- 
- Pessimistic upper bounds

## A Fundamental Tool: the Poincaré-Pontryagin Theorem



$$\begin{cases} \dot{x} = -\partial_y H(x, y) + \varepsilon f(x, y) \\ \dot{y} = \partial_x H(x, y) + \varepsilon g(x, y) \end{cases}$$

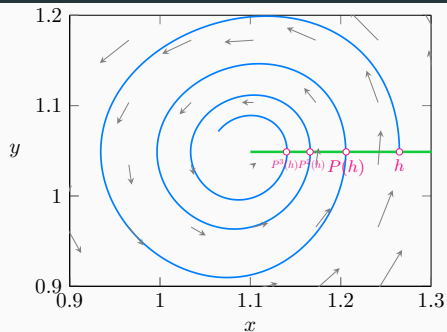
## A Fundamental Tool: the Poincaré-Pontryagin Theorem



- Poincaré first return map  $P(h)$

$$\begin{cases} \dot{x} = -\partial_y H(x, y) + \varepsilon f(x, y) \\ \dot{y} = \partial_x H(x, y) + \varepsilon g(x, y) \end{cases}$$

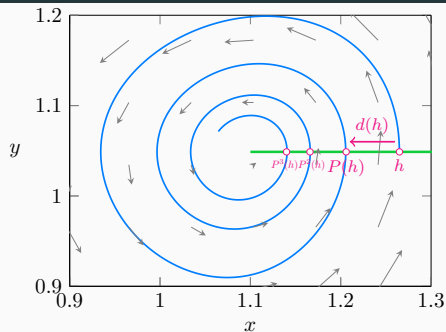
## A Fundamental Tool: the Poincaré-Pontryagin Theorem



- Poincaré first return map  $P(h)$

$$\begin{cases} \dot{x} = -\partial_y H(x, y) + \varepsilon f(x, y) \\ \dot{y} = \partial_x H(x, y) + \varepsilon g(x, y) \end{cases}$$

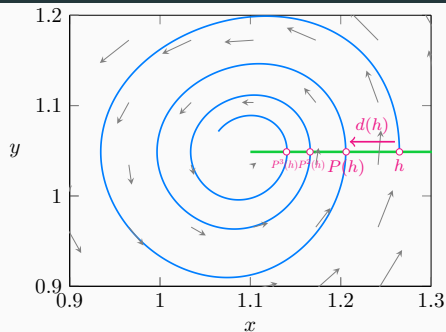
## A Fundamental Tool: the Poincaré-Pontryagin Theorem



- Poincaré first return map  $P(h)$
- Displacement  $d(h) = P(h) - h$

$$\begin{cases} \dot{x} = -\partial_y H(x, y) + \varepsilon f(x, y) \\ \dot{y} = \partial_x H(x, y) + \varepsilon g(x, y) \end{cases}$$

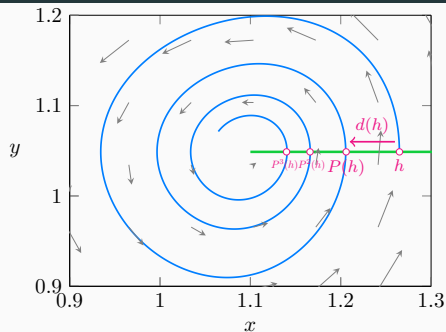
## A Fundamental Tool: the Poincaré-Pontryagin Theorem



- Poincaré first return map  $P(h)$
- Displacement  $d(h) = P(h) - h$
- Limit cycle  $\Leftrightarrow$  isolated zero of  $d$

$$\begin{cases} \dot{x} = -\partial_y H(x, y) + \varepsilon f(x, y) \\ \dot{y} = \partial_x H(x, y) + \varepsilon g(x, y) \end{cases}$$

## A Fundamental Tool: the Poincaré-Pontryagin Theorem



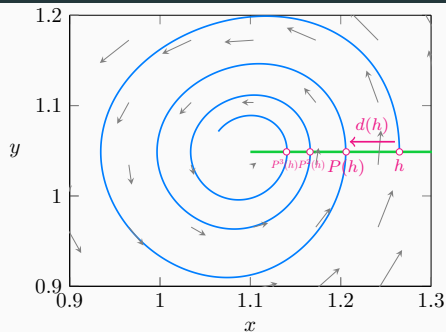
- Poincaré first return map  $P(h)$
- Displacement  $d(h) = P(h) - h$
- Limit cycle  $\Leftrightarrow$  isolated zero of  $d$
- Abelian integral  $\mathcal{I}(h)$ :

$$\oint_{H^{-1}(h)} f(x, y)dy - g(x, y)dx$$

$$\begin{cases} \dot{x} = -\partial_y H(x, y) + \varepsilon f(x, y) \\ \dot{y} = \partial_x H(x, y) + \varepsilon g(x, y) \end{cases}$$



# A Fundamental Tool: the Poincaré-Pontryagin Theorem



- Poincaré first return map  $P(h)$
- Displacement  $d(h) = P(h) - h$
- Limit cycle  $\Leftrightarrow$  isolated zero of  $d$
- Abelian integral  $\mathcal{I}(h)$ :

$$\oint_{H^{-1}(h)} f(x, y)dy - g(x, y)dx$$

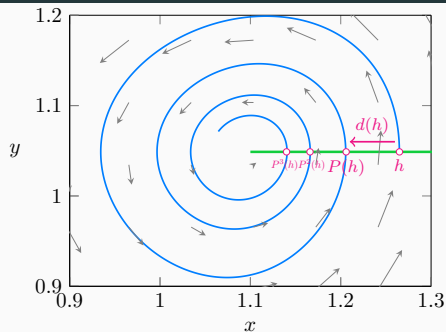
## Poincaré-Pontryagin theorem

The Abelian integral  $\mathcal{I}(h)$  approximates the displacement function  $d(h)$  for small  $\varepsilon$ :

$$d(h) = \varepsilon(\mathcal{I}(h) + O(\varepsilon)) \quad \text{when } \varepsilon \rightarrow 0$$

$$\begin{cases} \dot{x} = -\partial_y H(x, y) + \varepsilon f(x, y) \\ \dot{y} = \partial_x H(x, y) + \varepsilon g(x, y) \end{cases}$$

# A Fundamental Tool: the Poincaré-Pontryagin Theorem



- Poincaré first return map  $P(h)$
- Displacement  $d(h) = P(h) - h$
- Limit cycle  $\Leftrightarrow$  isolated zero of  $d$
- Abelian integral  $\mathcal{I}(h)$ :

$$\oint_{H^{-1}(h)} f(x, y)dy - g(x, y)dx$$

## Poincaré-Pontryagin theorem

The Abelian integral  $\mathcal{I}(h)$  approximates the displacement function  $d(h)$  for small  $\varepsilon$ :

$$d(h) = \varepsilon(\mathcal{I}(h) + O(\varepsilon)) \quad \text{when } \varepsilon \rightarrow 0$$

$$\begin{cases} \dot{x} = -\partial_y H(x, y) + \varepsilon f(x, y) \\ \dot{y} = \partial_x H(x, y) + \varepsilon g(x, y) \end{cases}$$

limit cycles  $\equiv$  changes of sign of  $\mathcal{I}(h)$   $\equiv$  simple zeros of  $\mathcal{I}(h)$