

Pade approximants to  $\exp(x)$

$$\exp(x) = \frac{N(x)}{D(x)} + O(x^{p+q+1}) \quad \begin{array}{l} N(x) \text{ of degree } p \\ D(x) \text{ poly. of degree } q \end{array}$$

Ex:  $p=q=2$

$$\exp(x) = \frac{ax^2+bx+c}{x^2+dx+e} + O(x^{p+q+1})$$



$$D(x)\exp(x) = N(x) + O(x^{p+q+1}) \text{ is linear in } a, b, c, d, e$$

$$(x^2+dx+e)\left(1+x+\frac{x^2}{2}+\frac{x^3}{3!}+\frac{x^4}{4!}+O(x^5)\right) - ax^2-bx-c = O(x^5)$$

$$\underbrace{e-c}_{=0} + \underbrace{(e+d-b)x}_{=0} + \underbrace{\left(1+d+\frac{e}{2}-a\right)x^2}_{=0} + \dots + \underbrace{\left(\frac{d}{6}+\frac{e}{24}\right)x^3}_{=0} + \underbrace{\left(\frac{d}{24}+\frac{e}{72}\right)x^4}_{=0} + O(x^5) = 0$$

Solution:  $N = x^2 + 6x + 12$        $\frac{N}{D} = \frac{x^2 + 6x + 12}{x^2 - 6x + 12}$   
 $D = x^2 - 6x + 12$

Backward error of Pade approximants:

$$f(x) = \log\left(\frac{e^{-x} N(x)}{D(x)}\right) = O(x^{p+q+1})$$

$$e^{f(x)} = e^{-x} \frac{N(x)}{D(x)} \Leftrightarrow e^{x+f(x)} = \frac{N(x)}{D(x)}$$

"backward error" setting:  $\frac{N(x)}{D(x)} = e^{\tilde{x}}$        $\tilde{x} = x + f(x)$

Let us apply everything to a matrix  $A$ :  $H = f(A)$

$Af = HA$  because  $f$  is a function of  $A$

$$e^H = e^{-A} N(A) D(A)^{-1} \Leftrightarrow e^{A+H} = e^A e^H = N(A) D(A)^{-1}$$

If  $\|H\| = O(u)\|A\|$  then  $N(A)D(A)^{-1}$  approximates  $e^A$  with an error that is comparable to intrinsic error (error induced by the approximation  $A \approx f(A)$ )       $u \approx 10^{-16}$

$$\tilde{A} \text{ for } (A) \quad \tilde{A}_{ij} = A_{ij}(1 + \epsilon_{ij}) \quad |\epsilon_{ij}| \leq u \approx 10^{-16}$$

$$\|\tilde{A} - A\| = \|A\| \cdot O(u)$$

$\Rightarrow$  the algorithm is "as good as it could get" on a machine (still does not ensure that it is accurate:  $e^{A+H}$  could be very far from  $e^A$ , if the problem is ill-conditioned)

$$f(x) = c_1 x^{p+q+1} + c_2 x^{p+q+2} + c_3 x^{p+q+3} + \dots \quad (\text{Taylor expansion})$$

$$H = f(A) = c_1 A^{p+q+1} + c_2 A^{p+q+2} + c_3 A^{p+q+3} + \dots$$

$$\|H\| = |c_1| \cdot \|A\|^{p+q+1} + |c_2| \cdot \|A\|^{p+q+2} + |c_3| \cdot \|A\|^{p+q+3} + \dots$$

If  $p=q=13$  and  $\|A\| \leq 5.4$ , then  $\frac{\|H\|}{\|A\|} \leq u \approx 10^{-16}$   
 [Higham book p.244]

for  $p=q$ ,  $N(x) = r(x^2) + xS(x^2)$  for some polynomials  $r, S$  of degree 6  
 $D(x) = r(x^2) - xS(x^2)$

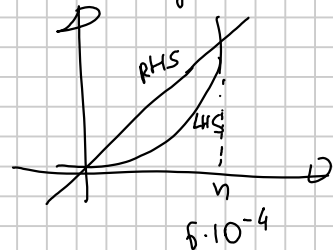
you can compute  $D/N$  with Paterson-Stockmeyer with 6 matrix multiplications

$$r(x^2) = (r_0 + r_1 x^2 + r_2 x^4 + r_3 x^6) + x^2 (r_4 x^2 + r_5 x^4 + r_6 x^6)$$

3 mults to form  $x^2, x^4, x^6$   
 1 for  $r(x^2)$       1 for  $xS(x^2)$   
 1 for  $S(x^2)$

$n = \|A\|$  we need to find  $n$  s.t. this inequality holds

$$\|H\| \leq |c_5| n^5 + |c_6| n^6 + \dots + |c_{30}| n^{30} \leq 10^{-16} \cdot n$$



$\Rightarrow$  for  $\|A\| \leq 6 \cdot 10^{-4}$ , the Padé approximant of degree (2,2) has backward error smaller than machine precision

What happens if  $\|A\| \geq 5.4$

Trick:  $e^x = (e^{\frac{1}{2}x})^2$        $\exp(A) = (\exp(\frac{1}{2}A))^2$

Algorithm: 1. Find  $k$  s.t.  $\|\frac{1}{2^k} A\| \leq 5.4$

"Scaling and squaring"

2. Compute  $\exp(\frac{1}{2^k} A) \approx N(\frac{1}{2^k} A) D(\frac{1}{2^k} A)^{-1} = F$   
 (Padé of degree  $p=q=13$ )

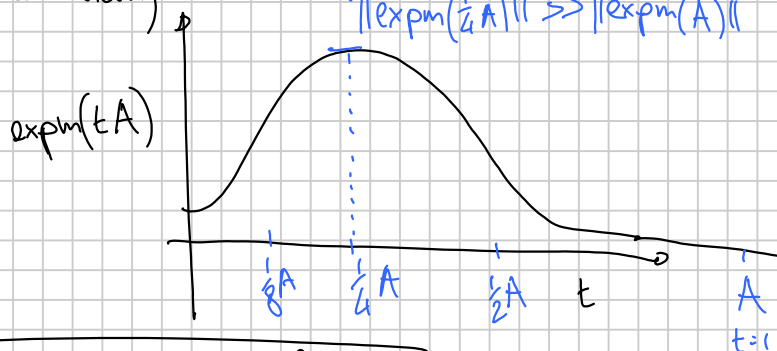
3.  $\exp(A) = F^{2^k}$  by repeated squaring

(This is what Matlab's `expm` does)

(If  $\|A\|$  is already much smaller than 5.4, it may use Padé expansions with degrees  $< 13$ )

(Other trick: if  $A$  is triangular,  $\expm(A)$  is triangular, too

On its diagonal, we should have  $\exp(A_{ii})$ ; we can compute them out and replace them.)



Conditioning of matrix functions

Condition number of  $y = f(x)$ :  $x \in \mathbb{R}^m$   
 $y \in \mathbb{R}^n$

$\Leftrightarrow$  how far away can  $\tilde{y} = f(\tilde{x})$  be from  $y$  if  $\tilde{x} - x$  is small?

$$\limsup_{\tilde{x} \rightarrow x} \frac{\|f(\tilde{x}) - f(x)\|}{\|\tilde{x} - x\|}$$



If  $f$  differentiable,  $f(\tilde{x}) = f(x) + J_{f,x} \cdot (\tilde{x} - x) + o(\|\tilde{x} - x\|)$

Absolute c.n.

$$\limsup_{\|\tilde{x} - x\| \rightarrow 0} \frac{\|J_{f,x}(\tilde{x} - x)\|}{\|\tilde{x} - x\|} + o(1) = \|J_{f,x}\|$$

(matrix norm)

Relative cond. number

$$K_{rel}(f,x) = \frac{\frac{\|f(\tilde{x}) - f(x)\|}{\|\tilde{x} - x\|}}{\frac{\|f(x)\|}{\|x\|}} = K_{abs}(f,x) \cdot \frac{\|x\|}{\|f(x)\|} = \|J_{f,x}\| \cdot \frac{\|x\|}{\|f(x)\|}$$

If  $f(A)$  is a matrix function, what is its Jacobian as a function  $\mathbb{R}^{n^2} \rightarrow \mathbb{R}^{n^2}$  ( $J_{f,A} \in \mathbb{R}^{n^2 \times n^2}$ )

Def: the Fréchet derivative: ↕

The Fréchet derivative of a matrix function  $f$  is the linear operator

$$L_{f,x}: \mathbb{R}^{n \times n} \rightarrow \mathbb{R}^{n \times n} \text{ such that}$$

$$f(X+E) = f(X) + \underbrace{L_{f,x}(E)} + o(\|E\|)$$

Example:  $f(x) = x^2 \rightarrow$  matrix function  $f(X) = X^2$

$$(X+E)^2 = \underbrace{X^2}_{f(X)} + \underbrace{XE+EX}_{\text{linear in } E} + \underbrace{E^2}_{o(\|E\|)}$$

The Fréchet derivative of  $f$  in  $X$  is the map  $L_{f,x}$  s.t.

$$L_{f,x}(E) = \underline{XE+EX}$$

(It is linear in  $E$ :  $L_{f,x}(E+F) = L_{f,x}(E) + L_{f,x}(F)$ )

$$L_{f,x}(\alpha E) = \alpha L_{f,x}(E)$$

$$\mathbb{R}^{n \times n} \cong \mathbb{R}^{n^2}$$

$\text{vec}(X) \in \mathbb{R}^{n^2}$  different point of view:

$\hat{L}$  is the Jacobian of  $\text{vec } f(\text{vec}^{-1}(x))$

$\hat{L}_{f,x}$  is the map  $\text{vec } E \rightarrow \text{vec } L_{f,x}(E) = \text{vec}(XE+EX)$

$$= \text{vec}(XE) + \text{vec}(EX)$$

$$= (I \otimes X) \text{vec } E + (X^T \otimes I) \text{vec } E$$

$$= (I \otimes X + X^T \otimes I) \text{vec } E$$

(recall that  $\text{vec}(AXB) = (B^T \otimes A) \text{vec}(X)$ )

$$\hat{L}_{f,x} = \underline{I \otimes X + X^T \otimes I} \in \mathbb{R}^{n^2 \times n^2}$$