

Advanced Topics from Scientific Computing  
TU Berlin Winter 2023/24  
Notebook 06

 Jürgen Fuhrmann

## Contents

### Forward mode automatic differentiation

- Dual numbers
- Dual numbers in Julia
- ForwardDiff.jl

### Solving nonlinear systems of equations

- Fixpoint iteration scheme
  - Definition of  $M(u)$
- Newton iteration scheme
  - Linear and quadratic convergence
  - newton1: Newton method with AD
  - dnewton: Damped Newton scheme
  - Parameter embedding
- NLsolve.jl
- Summary

# Forward mode automatic differentiation

---

## Dual numbers

---

The field of complex numbers  $\mathbb{C}$  extends the field of real numbers  $\mathbb{R}$  by introducing a number  $i$  with  $i^2 = -1$ .

Dual numbers are defined by extending the real numbers by formally introducing a number  $\varepsilon$  with  $\varepsilon^2 = 0$ :

$$\mathbb{D} = \{a + b\varepsilon \mid a, b \in \mathbb{R}\} = \left\{ \begin{pmatrix} a & b \\ 0 & a \end{pmatrix} \mid a, b \in \mathbb{R} \right\} \subset \mathbb{R}^{2 \times 2}$$

In contrast to real and complex numbers, dual numbers form a ring, not a field.

Evaluating polynomials on dual numbers: Let  $p(x) = \sum_{i=0}^n p_i x^i$ . Then

$$\begin{aligned} p(a + b\varepsilon) &= \sum_{i=0}^n p_i a^i + \sum_{i=1}^n i p_i a^{i-1} b \varepsilon \\ &= p(a) + b p'(a) \varepsilon \end{aligned}$$

- This can be generalized to any analytical function.  $\Rightarrow$  automatic evaluation of function and derivative at once
- $\Rightarrow$  *forward mode automatic differentiation*
- Multivariate dual numbers: generalization for partial derivatives

## Dual numbers in Julia

[Nathan Krislock](#) provided a simple dual number arithmetic example in Julia.

- Define a struct parametrized with type T. This is akin a template class in C++
- The type shall work with all methods working with `Number`
- In order to construct a Dual number from arguments of different types, allow promotion aka "parameter type homogenization"

```
1 begin
2     struct DualNumber{T} <: Number where {T <: Real}
3         value::T
4         deriv::T
5     end
6     DualNumber(v,d) = DualNumber(promote(v,d)...)
7 end;
```

Define a way to convert a `Real` to `DualNumber`

```
1 Base.promote_rule(::Type{DualNumber{T}}, ::Type{<:Real}) where T<:Real =
  DualNumber{T}
```

```
1 Base.convert(::Type{DualNumber{T}}, x::Real) where T<:Real = DualNumber(x,zero(T))
2
```

Constructing a dual number:

```
d = DualNumber(5, 4)
```

```
1 d=DualNumber(5,4)
```

Accessing its components:

```
(5, 4)
```

```
1 d.value,d.deriv
```

Simple arithmetic for dual numbers:

Add methods to the functions +, /, \*, -, inv which allow them to work for DualNumber

```
1 begin
2   import Base: +, /, *, -, inv
3   +(x::DualNumber, y::DualNumber) =
4       DualNumber(x.value + y.value, x.deriv + y.deriv)
5
6   -(y::DualNumber) = DualNumber(-y.value, -y.deriv)
7
8   -(x::DualNumber, y::DualNumber) = x + -y
9
10  *(x::DualNumber, y::DualNumber) =
11      DualNumber(x.value*y.value, x.value*y.deriv + x.deriv*y.value)
12
13  inv(y::DualNumber{T}) where T<:Union{Integer, Rational} =
14      DualNumber(1/y.value, (-y.deriv)//y.value^2)
15
16  inv(y::DualNumber{T}) where T<:Union{AbstractFloat,AbstractIrrational} =
17      DualNumber(1/y.value, (-y.deriv)/y.value^2)
18
19  /(x::DualNumber, y::DualNumber) = x*inv(y)
20 end;
21
```

```
1 Base.sin(x::DualNumber{T}) where T= DualNumber(sin(x.value),cos(x.value)*x.deriv);
```

```
1 Base.log(x::DualNumber{T}) where T = DualNumber(log(x.value),x.deriv/x.value)
```

Define a function for comparison with known derivative:

testdual (generic function with 1 method)

```
1 function testdual(x,f,df)
2   xdual=DualNumber(x,1)
3   fdual=f(xdual)
4   _f=f(x)
5   _df=df(x)
6   err=_df-fdual.deriv
7   (f=_f,f_dual=fdual.value),(df=_df,df_dual=fdual.deriv), (error=err,)
8 end
```

Polynomial expressions:

p (generic function with 1 method)

```
1 p(x)=x^3+2x+1
```

dp (generic function with 1 method)

```
1 dp(x)=3x^2+2
```

```
((f = 34, f_dual = 34), (df = 29, df_dual = 29), (error = 0))
```

```
1 testdual(3,p,dp)
```

Standard functions:

```
((f = 0.420167, f_dual = 0.420167), (df = 0.907447, df_dual = 0.907447), (error = 0.0))
```

```
1 testdual(13,sin,cos)
```

```
((f = 2.56495, f_dual = 2.56495), (df = 0.0769231, df_dual = 0.0769231), (error = 0.0))
```

```
1 testdual(13,log, x->1/x)
```

Function composition:

```
((f = -0.506366, f_dual = -0.506366), (df = 17.2464, df_dual = 17.2464), (error = 0.0))
```

```
1 testdual(10,x->sin(x^2),x->2x*cos(x^2))
```

If we apply dual numbers in the right way, we can do calculations with derivatives of complicated nonlinear expressions without the need to write code to calculate derivatives.

## ForwardDiff.jl

```
1 using ForwardDiff, LinearAlgebra
```

The [ForwardDiff.jl](#) package provides a full implementation of these facilities.

testdual1 (generic function with 1 method)

```
1 function testdual1(x,f,df)
2     _df=df(x)
3     _df_dual=ForwardDiff.derivative(f,x)
4     (f=f(x),df=_df,df_dual=_df_dual, error=abs(_df-_df_dual))
5 end
```

```
(f = 0.14112, df = -0.989992, df_dual = -0.989992, error = 0.0)
```

```
1 testdual1(3,sin,cos)
```

Let us plot some complicated function:

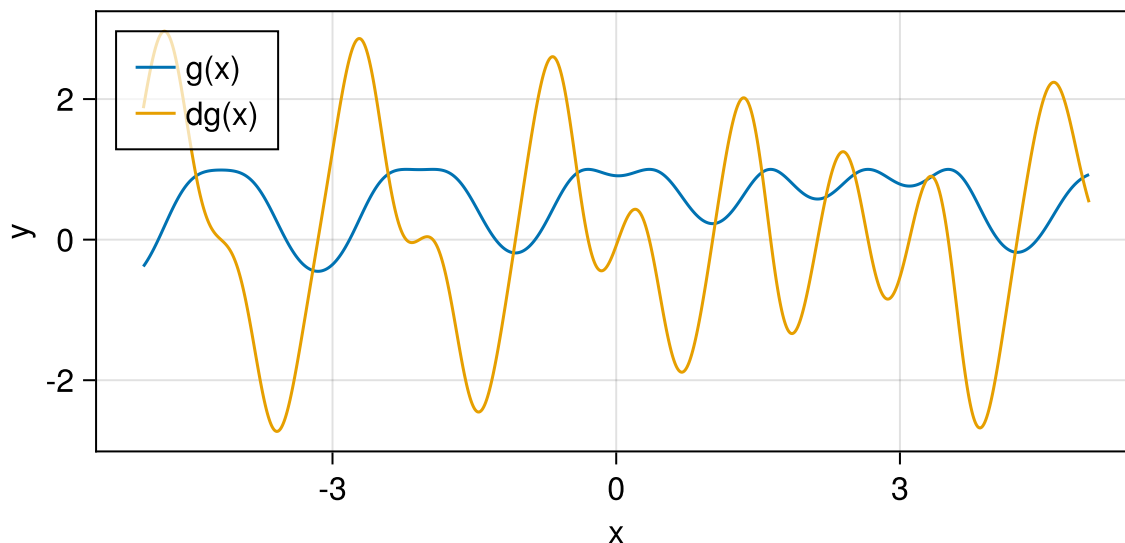
g (generic function with 1 method)

```
1 g(x)=sin(exp(0.2*x))+cos(3x)
```

dg (generic function with 1 method)

```
1 dg(x)=ForwardDiff.derivative(g,x)
```

```
1 begin
2     using CairoMakie, Colors
3     CairoMakie.activate!(type="svg")
4 end
```



```
1 let
2     X=(-5:0.01:5)
3     fig=Figure(resolution=(600,300))
4     axis=Axis(fig[1,1],xlabel="x",ylabel="y")
5     lines!(axis,X,g.(X),label="g(x)")
6     lines!(axis,X,dg.(X),label="dg(x)")
7     axislegend(axis,bgcolor = RGBA(1.0, 1.0, 1.0, 0.7),position=:lt)
8     fig
9 end
```

## Solving nonlinear systems of equations

---

Let  $A_1 \dots A_n$  be functions depending on  $n$  unknowns  $u_1 \dots u_n$ . Solve the system of nonlinear equations:

$$A(u) = \begin{pmatrix} A_1(u_1 \dots u_n) \\ A_2(u_1 \dots u_n) \\ \vdots \\ A_n(u_1 \dots u_n) \end{pmatrix} = \begin{pmatrix} f_1 \\ f_2 \\ \vdots \\ f_n \end{pmatrix} = f$$

$A(u)$  can be seen as a nonlinear operator  $A : D \rightarrow \mathbb{R}^n$  where  $D \subset \mathbb{R}^n$  is its domain of definition.

There is no analogon to Gaussian elimination, so we need to solve iteratively.

## Fixpoint iteration scheme

Assume  $A(u) = M(u)u$  where for each  $u$ ,  $M(u) : \mathbb{R}^n \rightarrow \mathbb{R}^n$  is a linear operator.

Then we can define the iteration scheme: choose an initial value  $u_0$  and at each iteration step, solve

$$M(u^i)u^{i+1} = f$$

Terminate if

$$\|A(u^i) - f\| < \varepsilon \quad (\text{residual based})$$

or

$$\|u_{i+1} - u_i\| < \varepsilon \quad (\text{update based}).$$

- Large domain of convergence
- Convergence may be slow
- Smooth coefficients not necessary

fixpoint! (generic function with 1 method)

```

1 function fixpoint!(u,M,f; imax=100, tol=1.0e-10)
2     history=Float64[]
3     for i=1:imax
4         res=norm(M(u)*u-f)
5         push!(history,res)
6         if res<tol
7             return u,history
8         end
9         u=M(u)\f
10    end
11    error("No convergence after $imax iterations")
12 end
13

```

# Definition of $M(u)$

M (generic function with 1 method)

```
1 function M(u)
2     [ 1+1.2*(u[1]^2+u[2]^2)  -(u[1]^2+u[2]^2);
3       -(u[1]^2+u[2]^2)  1+1*(u[1]^2+u[2]^2)]
4 end
```

F = [1, 3]

```
1 F=[1,3]
```

[1.28822, 1.61348], [3.16228, 26.9072, 1.45019, 1.87735, 0.614397, 0.471544, 0.229973, 0.139343, 0.091045, 0.060485]

```
1 fixpt_result,fixpt_history=fixpoint!([0,0],M,F,imax=1000,tol=1.0e-10)
```

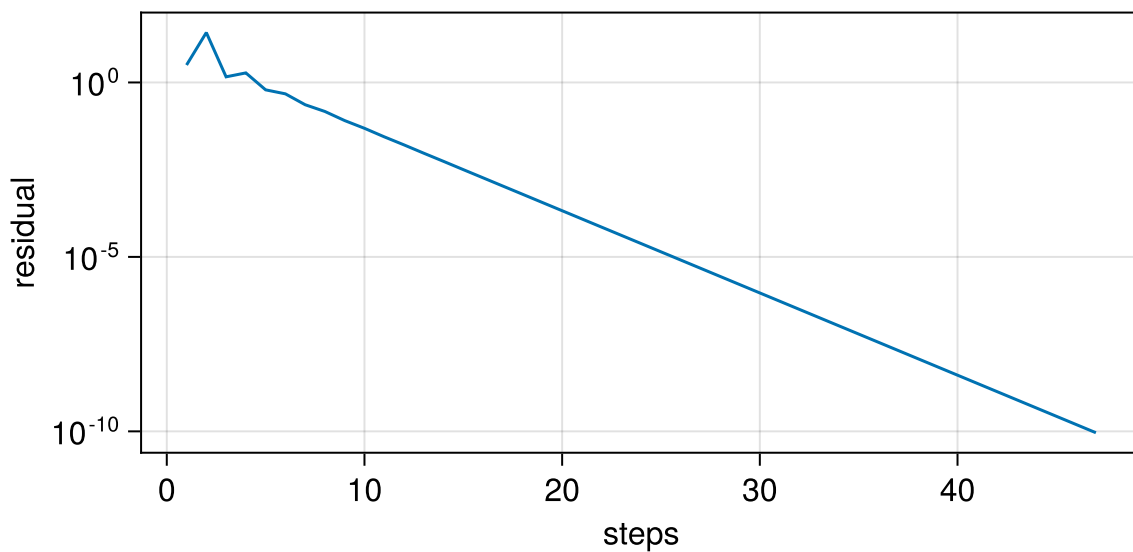
contraction (generic function with 1 method)

```
1 contraction(h)=h[2:end]./h[1:end-1]
```

```
1 function plothistory(history::Vector{<:Number})
2     fig=Figure(resolution=(600,300))
3     axis=Axis(fig[1,1],xlabel="steps",ylabel="residual",yscale=log10)
4     lines!(axis,1:length(history),history)
5     # axislegend(axis,bgcolor = RGBA(1.0, 1.0, 1.0, 0.7),position=:lt)
6     fig
7 end;
```

[8.50882, 0.0538958, 1.29456, 0.327268, 0.76749, 0.487702, 0.640077, 0.548586, 0.60068, 0.614397, 0.471544, 0.229973, 0.139343, 0.091045, 0.060485]

```
1 contraction(fixpt_history)
```



```
1 plothistory(fixpt_history)
```

[1.85807e-11, -8.93863e-11]

1 `M(fixpt_result)*fixpt_result-F`

## Newton iteration scheme

The fixed point iteration scheme assumes a particular structure of the nonlinear system. In addition, one would need to investigate convergence conditions for each particular operator. Can we do better?

Let  $A'(u)$  be the *Jacobi matrix* of first partial derivatives of  $A$  at point  $u$ :

$$A'(u) = (a_{kl})$$

'with

$$a_{kl} = \frac{\partial}{\partial u_l} A_k(u_1 \dots u_n)$$

Then, one calculates in the  $i$ -th iteration step:

$$u_{i+1} = u_i - (A'(u_i))^{-1}(A(u_i) - f)$$

One can split this as follows:

- Calculate residual:  $r_i = A(u_i) - f$
- Solve linear system for update:  $A'(u_i)h_i = r_i$
- Update solution:  $u_{i+1} = u_i - h_i$

General properties are:

- Potentially small domain of convergence - one needs a good initial value
- Possibly slow initial convergence
- Quadratic convergence close to the solution



# Linear and quadratic convergence

Let  $e_i = u_i - \hat{u}$ .

- Linear convergence: observed for e.g. linear systems: Asymptotically constant error contraction rate

$$\frac{\|e_{i+1}\|}{\|e_i\|} \sim \rho < 1$$

- Quadratic convergence:  $\exists i_0 > 0$  such that  $\forall i > i_0, \frac{\|e_{i+1}\|}{\|e_i\|^2} \leq M < 1$ .
  - As  $\|e_i\|$  decreases, the contraction rate decreases:

$$\frac{\frac{\|e_{i+1}\|}{\|e_i\|}}{\frac{\|e_i\|}{\|e_{i-1}\|}} = \frac{\|e_{i+1}\|}{\frac{\|e_i\|^2}{\|e_{i-1}\|}} \leq \|e_{i-1}\| M$$

- In practice, we can watch  $\|r_i\|$  or  $\|h_i\|$

## newton1: Newton method with AD

This is the situation where we could apply automatic differentiation for vector functions of vectors.

A1 (generic function with 1 method)

```
1 A1(u)=M(u)*u
```

newton1 (generic function with 1 method)

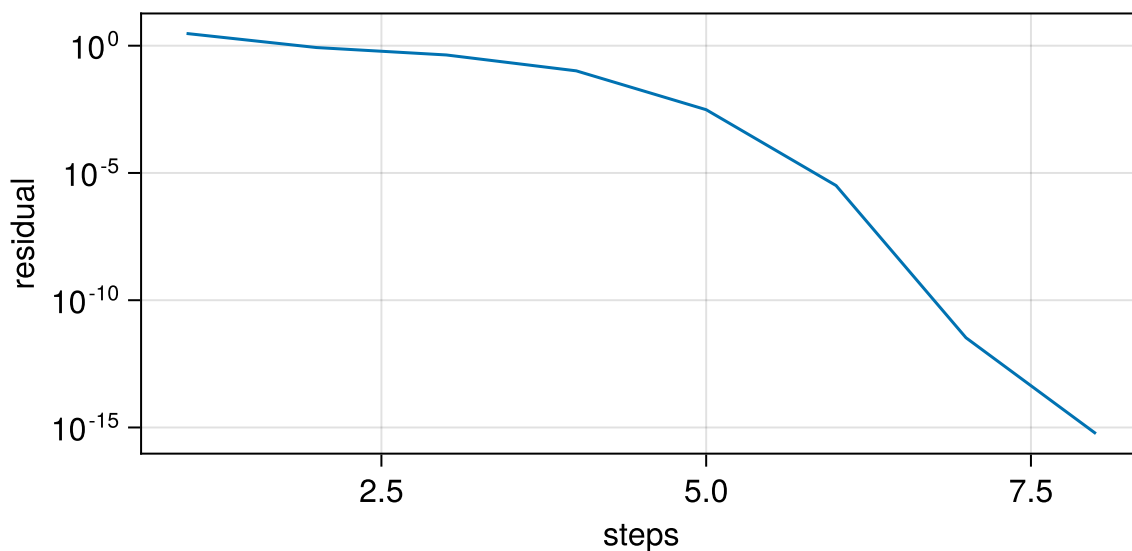
```

1 function newton1(A,b,u0; tol=1.0e-12, maxit=100)
2     history=Float64[]
3     u=copy(u0)
4     it=0
5     converged=false
6     while !converged && it<maxit
7         res=A(u)-b
8         jac=ForwardDiff.jacobian((v)->A(v)-b ,u)
9         h=jac\res
10        u-=h
11        nm=norm(h)
12        push!(history,nm)
13        it=it+1
14        if nm<tol
15            converged=true
16        end
17    end
18    if converged
19        return u,history
20    else
21        throw("convergence failed")
22    end
23 end

```

([1.28822, 1.61348], [3.02185, 0.846373, 0.432681, 0.102853, 0.0030576, 3.19945e-6, 3.35:

```
1 newton_result1,newton_history1=newton1(A1,F,[0,0.1],tol=1.e-13)
```



```
1 plotohistory(newton_history1)
```

Calculate function and derivative at once ?

Let us take a more complicated example with an operator dependent on a parameter  $\lambda$  which allows to adjust the "severity" of the nonlinearity. For  $\lambda=0$ , it is linear, for  $\lambda=1$  it is strongly nonlinear.

A2λ (generic function with 1 method)

```
1 A2λ(x,λ)= [x[1]+10λ*x[1]^5+3λ*x[2]*x[3],
2           0.1*x[2]+10λ*x[2]^5-3λ*x[1]-x[3],
3           10λ*x[3]^5+10λ*x[1]*x[2]*x[3]+x[3]/100]
```

A2 (generic function with 1 method)

```
1 A2(x)=A2λ(x,1)
```

F2 = [0.1, 0.1, 0.1]

```
1 F2=[0.1,0.1,0.1]
```

U02 = [1.0, 1.0, 1.0]

```
1 U02=[1,1.0,1.0]
```

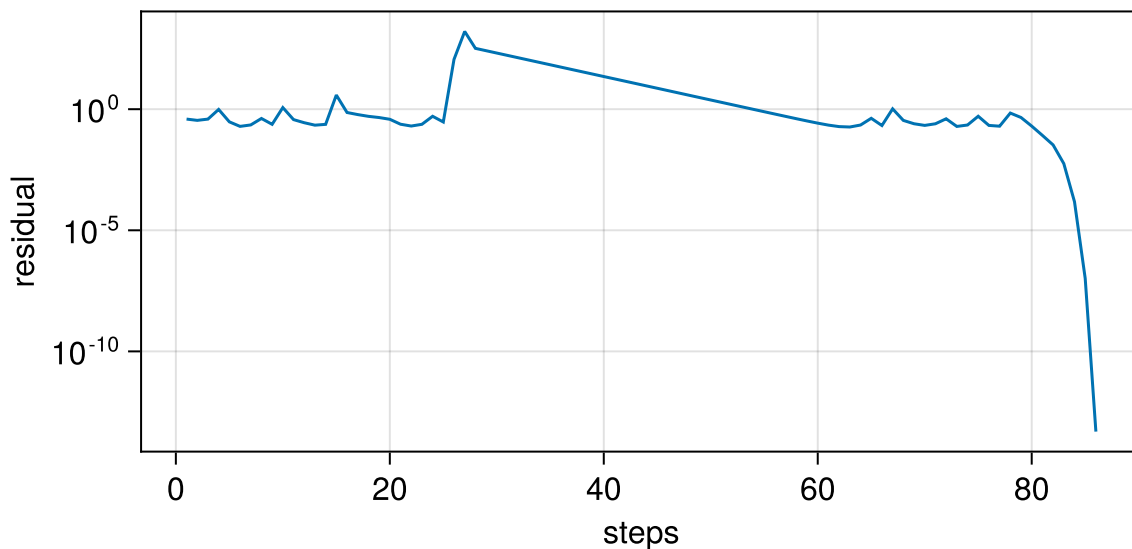
([-0.188484, 0.198519, 0.488388], [0.39077, 0.345694, 0.389908, 0.977557, 0.300465, 0.198519])

```
1 res2,hist2=newton1(A2,F2,U02)
```

[-2.77556e-17, -2.77556e-17, 0.0]

```
1 A2(res2)-F2
```

Newton steps: 86



```
1 plotohistory(hist2)
```

Here, we observe that we have to use lots of iteration steps and see a rather erratic behaviour of the residual. After  $\approx 80$  steps we arrive in the quadratic convergence region where convergence is fast.

# dnewton: Damped Newton scheme

There are many ways to improve the convergence behaviour and/or to increase the convergence radius in such a case. The simplest ones are:

- find a good estimate of the initial value
- damping: do not use the full update, but damp it by some factor which we increase during the iteration process until it reaches 1

dnewton (generic function with 1 method)

```

1 function dnewton(A,b,u0; tol=1.0e-12,maxit=100,damp=1,damp_growth=1)
2     result=DiffResults.JacobianResult(u0)
3     history=Float64[]
4     u=copy(u0)
5     it=1
6     while it<maxit
7         ForwardDiff.jacobian!(result,(v)->A(v)-b ,u)
8         res=DiffResults.value(result)
9         jac=DiffResults.jacobian(result)
10        h=jac\res
11        u.-=damp*h
12        nm=norm(h)
13        push!(history,nm)
14        if nm<tol
15            return u,history
16        end
17
18        it=it+1
19        damp=min(damp*damp_growth,1.0)
20    end
21    throw("convergence failed")
22 end

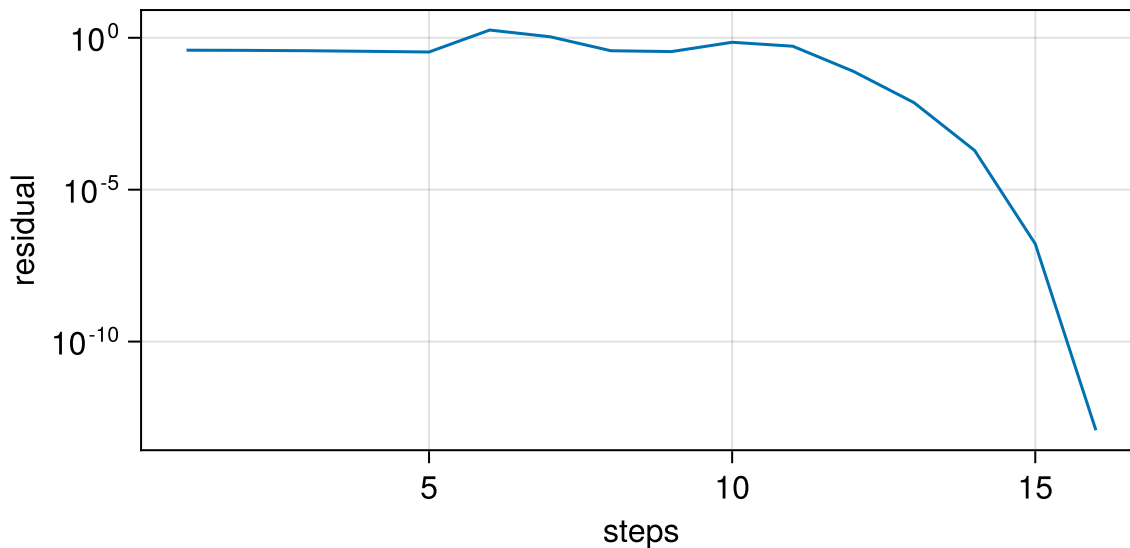
```

In this implementation, we also try to save work by evaluating result and Jacobian once.

```
([-0.188484, 0.198519, 0.488388], [0.39077, 0.38541, 0.375394, 0.358292, 0.340649, 1.798
```

```
1 res3,hist3=dnewton(A2,F2,U02,damp=0.1,damp_growth=2,maxit=1000)
```

Newton steps: 16



```
1 plothistory(hist3)
```

```
[-2.77556e-17, -2.77556e-17, 0.0]
```

```
1 A2(res3)-F2
```

The example shows: damping indeed helps to improve the convergence behaviour. If we would keep the damping parameter less than 1, we lose the quadratic convergence behavior.

A more sophisticated strategy would be line search: automatic detection of a damping factor which prevents the residual from increasing.

## Parameter embedding

Another option is the use of parameter embedding for parameter dependent problems.

- Problem: solve  $A(u_\lambda, \lambda) = f$  for  $\lambda = 1$ .
- Assume  $A(u_0, 0)$  can be easily solved.
- Choose step size  $\delta$

1. Solve  $A(u_0, 0) = f$
2. Set  $\lambda = 0$
3. Solve  $A(u_{\lambda+\delta}, \lambda + \delta) = f$  with initial value  $u_\lambda$
4. Set  $\lambda = \lambda + \delta$
5. If  $\lambda < 1$  repeat with 3.

- If  $\delta$  is small enough, we can ensure that  $u_\lambda$  is a good initial value for  $u_{\lambda+\delta}$ .
- Possibility to adapt  $\delta$  depending on Newton convergence



```
1 using NLSolve
```

```
nlres1 = Results of Nonlinear Solver Algorithm
* Algorithm: Trust-region with dogleg and autoscaling
* Starting Point: [1.0, 1.0, 1.0]
* Zero: [0.05758245053199316, 0.48399543292760894, 0.04126490097524311]
* Inf-norm of residuals: 0.088086
* Iterations: 1000
* Convergence: false
*  $|x - x'| < 0.0e+00$ : false
*  $|f(x)| < 1.0e-08$ : false
* Function Calls (f): 79
* Jacobian Calls (df/dx): 38
```

```
1 nlres1=nlsolve(u->A2λ(u,1.0)-F2, U02)
```

```
[0.0175049, -2.60122e-5, -0.0880858]
```

```
1 A2λ(nlres1.zero,1.0)-F2
```

```
nlres2 = Results of Nonlinear Solver Algorithm
* Algorithm: Newton with line-search
* Starting Point: [1.0, 1.0, 1.0]
* Zero: [-0.18848435786947373, 0.198519144942218, 0.4883882611017444]
* Inf-norm of residuals: 0.000000
* Iterations: 239
* Convergence: true
*  $|x - x'| < 0.0e+00$ : false
*  $|f(x)| < 1.0e-08$ : true
* Function Calls (f): 240
* Jacobian Calls (df/dx): 240
```

```
1 nlres2=nlsolve(u->A2λ(u,1.0)-F2, U02, method=:newton)
```

```
[-1.12965e-14, 8.32667e-17, 7.83734e-13]
```

```
1 A2λ(nlres2.zero,1.0)-F2
```

```
nlres3 = Results of Nonlinear Solver Algorithm
* Algorithm: Newton with line-search
* Starting Point: [1.0, 1.0, 1.0]
* Zero: [-0.18848435786937287, 0.19851914494226677, 0.48838826110144995]
* Inf-norm of residuals: 0.000000
* Iterations: 85
* Convergence: true
*  $|x - x'| < 0.0e+00$ : false
*  $|f(x)| < 1.0e-08$ : true
* Function Calls (f): 86
* Jacobian Calls (df/dx): 86
```

```
1 nlres3=nlsolve(u->A2λ(u,1.0)-F2, U02, method=:newton, autodiff=:forward)
```

```
[-7.91034e-15, 5.27356e-16, 1.06304e-13]
```

```
1 A2λ(nlres3.zero,1.0)-F2
```

# Summary

---

- Newton method with increasing damping + update based convergence control is rather robust
    - I use this in my everyday work
  - Additional parameter embedding can help to solve even strongly nonlinear problems
  - NLSolve.jl provides a convenient default first stop for solving nonlinear systems in Julia, it relies on a number of peer reviewed strategies
- 

```
1 using HypertextLiteral, PlutoUI
```