# Automatic differentiation

**Advanced Topics from Scientific Computing**
**TU Berlin Winter 2024/25**
**Notebook 06**
(cc) BY-SA **Jürgen Fuhrmann**

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

## Contents

Automatic differentiation (AD) has become more and more widely used. It is one of main ingredients of machine learning. This lecture tries to explain some basic ideas on how it works.

For further reading: the Wikipedia article on AD is well resourced.

# Reverse mode AD

.................................................................................................

Let us define a function and see how we can calculate it derivative:

f (generic function with 2 methods)

```
1  f(x, y) = x^2 + 2x * y + 3y
```

```
1  x = 10.0; y = 11.0;
```

353.0
```
1  f(x, y)
```

fdf (generic function with 1 method)

```
1  function fdf(x, y)
2      #! format: off
3      v1 = x^2;    ∂x_v1 = 2x;    ∂y_v1 = 0
4      v2 = 2x * y; ∂x_v2 = 2y;    ∂y_v2 = 2x
5      v3 = 3y;     ∂x_v3 = 0;     ∂y_v3 = 3
6
7      v4 = v1 + v2 + v3
8      ∂x_v4 = ∂x_v1 + ∂x_v2 + ∂x_v3
9      ∂y_v4 = ∂y_v1 + ∂y_v2 + ∂y_v3
10
11     v4, ∂x_v4, ∂y_v4
12     #! format: on
13 end
```

(353.0, 42.0, 23.0)
```
1  fdf(x, y)
```

In order to caculate the derivatives, one essentially parses the expression and recursively applies differentiation rules.

# A simple implementation of reverse AD

Here, we provide a simple example how this can be implemented. We modify the implementation by Simeon Schaub in his ReverseModePluto example.

```
1 begin
2     using AbstractTrees
3     using LinearAlgebra
4 end
```

## Tracked number type

We use the Julia multiple dispatch system in order to design a number type (which we call `Tracked`) which allows to track operations and apply derivative rules.

```
1 begin
2     struct Tracked{T} <: Number
3         val::T    # value
4         name::Symbol # label just for printing
5         df::Union{Function, Nothing} # derivative rule combining deps
6         deps::Vector{Tracked} # Tracked values of dependent expressions.
7     end
8
9     # Construct a tracked value from with optional name from some number
10    function Tracked{T}(x, name = gensym()) where {T}
11        Tracked{T}(x, name, nothing, Tracked[])
12    end
13
14    # Constructor which automatically takes type from number
15    Tracked(x::T, name = gensym()) where {T} = Tracked{T}(x, name)
16 end;
```

## Equality, conversion, promotion

When are two `Tracked` instances equal ?

```
1 Base.:(==)(x::Tracked, y::Tracked) = x === y
```

Lift some rules for handling number types to `Tracked`:

Allow to convert tracked values with different type parameters

```
1 function Base.convert(T::Type{Tracked{S}}, x::Tracked) where {S}
2     T(convert(S, x.val), x.name, x.df, x.deps)
3 end
```

Promote rule for handling combinations of tracked values with different type parameters

```
1 function Base.promote_rule(::Type{Tracked{S}},
2                            ::Type{T}) where {S <: Number, T <: Number}
3     Tracked{promote_type(S, T)}
4 end
```

## Printing tracked numbers as tree

```
1  AbstractTrees.children(x::Tracked) = x.deps
```

```
1  function AbstractTrees.printnode(io::IO, node::Tracked)
2      if isnothing(node.df)
3          if Base.isgensym(node.name)
4              print(io, " ← $(node.val)")
5          else
6              print(io, "[$(node.name)] ← $(node.val)")
7          end
8      else
9          print(io, "[$(node.name)] = $(node.val)")
10     end
11 end
```

```
1  Base.show(io::IO, ::MIME"text/plain", x::Tracked) = print_tree(io, x)
```

## Arithmetics of tracked numbers

Now we can print tracked numbers and define operations using mutiple dispatch for the basic operations on numbers.

```
tx = [x] ← 10.0
```
```
1  tx = Tracked(x, :x)
```

```
ty = [y] ← 11.0
```
```
1  ty = Tracked(y, :y)
```

```
1  function Base.:+(x::Tracked, y::Tracked)
2      Tracked(x.val + y.val, :+, Δ -> (Δ, Δ), Tracked[x, y])
3  end
```

Here, we see the tree for addition

```
[+] = 21.0
├─ [x] ← 10.0
└─ [y] ← 11.0
```
```
1  tx + ty
```

```
1  function Base.:-(x::Tracked, y::Tracked)
2      Tracked(x.val - y.val, :-, Δ -> (Δ, -Δ), Tracked[x, y])
3  end
```

```
1  function Base.:*(x::Tracked, y::Tracked)
2      Tracked(x.val * y.val, :*, Δ -> (Δ * y.val, x.val * Δ), Tracked[x, y])
3  end
```

```
1  function Base.:^(x::Tracked, n::Int)
2      Tracked(x.val^n, Symbol("^$n"), Δ -> (Δ * n * x.val^(n - 1),), Tracked[x])
3  end
```

```
1  function Base.:/(x::Tracked, y::Tracked)
2      Tracked(x.val / y.val, :/, Δ -> (Δ / y.val, -Δ * x.val / y.val^2),
3              Tracked[x, y])
4  end
```

These operations are sufficient to be applied to our initial function:

```
[+] = 353.0
├─ [+] = 320.0
│  ├─ [^2] = 100.0
│  │  └─ [x] ← 10.0
│  └─ [*] = 220.0
│     ├─ [*] = 20.0
│     │  ├─ ← 2.0
│     │  └─ [x] ← 10.0
│     └─ [y] ← 11.0
└─ [*] = 33.0
   ├─ ← 3.0
   └─ [y] ← 11.0
```

```
1  f(tx, ty)
```

# Collecting the gradient

Set up a dictionary where tracked values are the keys:

```
1  function grad(f::Tracked)
2      d = Dict{Any, Any}(f => 1)
3      # Depth first tree search (DFS)
4      # recursively traverses all dependents, parents before children:
5      for x in PreOrderDFS(f)
6          # ignore untracked variables like constants
7          if x.df === nothing
8              continue
9          end
10         dy = x.df(d[x]) # Evaluate "pullback": parents have been visited before,
11         # so they have dictionary entries
12         for (yᵢ, dyᵢ) in zip(x.deps, dy)
13             # store the gradient in d
14             # if we have already stored a gradient for
15             # this variable, we need to add them
16             if haskey(d, yᵢ)
17                 d[yᵢ] = d[yᵢ] + dyᵢ
18             else
19                 d[yᵢ] = dyᵢ
20             end
21         end
22     end
23     return d
24  end;
```

```
1  grad(f::Tracked, x::Tracked) = grad(f)[x];
```

```
tf = [+] = 353.0
     ├── [+] = 320.0
     │    ├── [^2] = 100.0
     │    │    └── [x] ← 10.0
     │    └── [*] = 220.0
     │         ├── [*] = 20.0
     │         │    ├── ← 2.0
     │         │    └── [x] ← 10.0
     │         └── [y] ← 11.0
     └── [*] = 33.0
          ├── ← 3.0
          └── [y] ← 11.0
```

```
1  tf = f(tx, ty)
```

```
42.0
```

```
1  grad(tf, tx)
```

```
23.0
```

```
1  grad(tf, ty)
```

```
1  function Base.sin(x::Tracked)
2      Tracked(sin(x.val), :sin, Δ -> (Δ * cos(x.val),), Tracked[x])
3  end
```

```
1  function Base.cos(x::Tracked)
2      Tracked(cos(x.val), :cos, Δ -> (-Δ * sin(x.val),), Tracked[x])
3  end
```

```
1  function Base.exp(x::Tracked)
2      Tracked(exp(x.val), :exp, Δ -> (Δ * exp(x.val),), Tracked[x])
3  end
```

So now, we can perform AD for functions:

```
ad (generic function with 1 method)
```

```
1  function ad(f, x)
2      tx = Tracked(x)
3      grad(f(tx), tx)
4  end
```

```
g (generic function with 1 method)
```
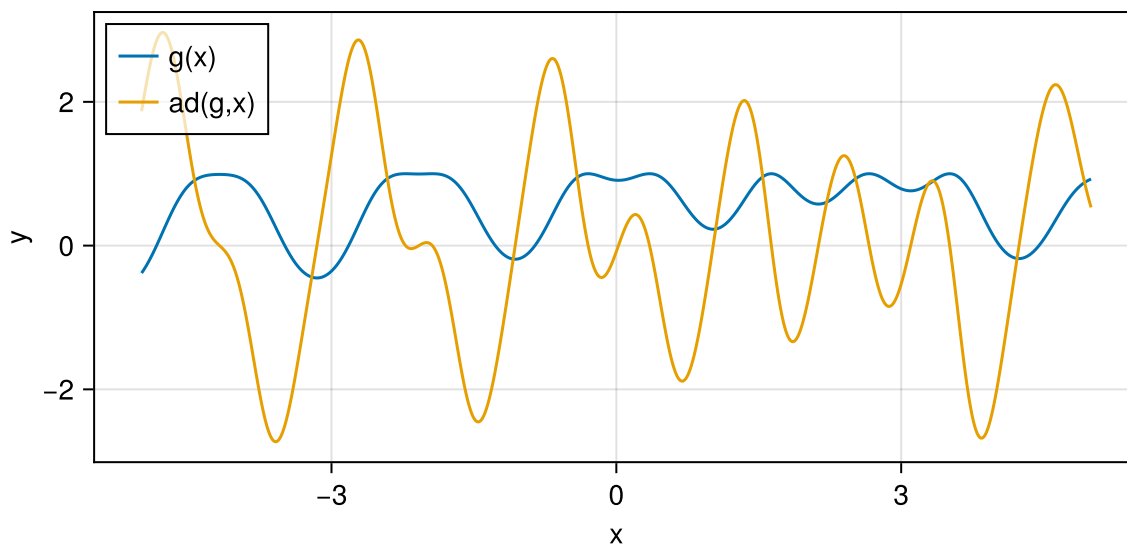
```
1  g(x) = sin(exp(0.2 * x) + cos(3x))
```

```
[sin] = 0.9521277180368776
 └── [+] = 7.543307548818235
      ├── [exp] = 7.38905609893065
      │    └── [*] = 2.0
      │         ├── ← 0.2
      │         └── [x] ← 10.0
      └── [cos] = 0.15425144988758405
           └── [*] = 30.0
                ├── ← 3.0
                └── [x] ← 10.0
```

```
1  g(tx)
```

```
1  let
2      X = (-5:0.01:5)
3      fig = Figure(; size = (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, ad.(g, X); label = "ad(g,x)")
7      axislegend(axis; backgroundcolor = RGBA(1.0, 1.0, 1.0, 0.7), position = :lt)
8      fig
9  end
```

In general, reverse mode AD implementations use something similar to our Tracked tree, which "records" the function and therefore often is calles a "Tape". This means that in this case AD computation requires to first record the tape, which may result in certain memory overhead.

# ReverseDiff.jl and DifferentiationInterface.jl

Julia provides an implementation of reverse mode AD in ReverseDiff.jl and a general interface to AD via DifferentiationInterface.jl

```
1  using ReverseDiff, DifferentiationInterface
```

The AD function needs a vector input, so we define another method for our initial `f`:

f (generic function with 1 method)
```
1  f(X::AbstractVector) = f(X[1], X[2])
```

(353.0, [42.0, 23.0])
```
1  value_and_gradient(f, AutoReverseDiff(), [x, y])
```

# Forward mode AD

We cover one particular formulation of forward mode AD based on so-called dual numbers.

# 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

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

- 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`

```julia
1  function Base.promote_rule(::Type{DualNumber{T}},
2                             ::Type{<:Real}) where {T <: Real}
3      DualNumber{T}
4  end
```

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

Constructing a dual number:

d =    DualNumber(5, 4)

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

Accessing its components:

(5, 4)

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

Simple arithmetic for dual numbers:

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

```julia
1  begin
2      import Base: +, /, *, -, inv
3      function +(x::DualNumber, y::DualNumber)
4          DualNumber(x.value + y.value, x.deriv + y.deriv)
5      end
6
7      -(y::DualNumber) = DualNumber(-y.value, -y.deriv)
8
9      -(x::DualNumber, y::DualNumber) = x + -y
10
11     function *(x::DualNumber, y::DualNumber)
12         DualNumber(x.value * y.value, x.value * y.deriv + x.deriv * y.value)
13     end
14
15     function inv(y::DualNumber{T}) where {T <: Union{Integer, Rational}}
16         DualNumber(1 // y.value, (-y.deriv) // y.value^2)
17     end
18
19     function inv(y::DualNumber{T}) where {T <: Union{AbstractFloat,
20                                                      AbstractIrrational}}
21         DualNumber(1 / y.value, (-y.deriv) / y.value^2)
22     end
23
24     /(x::DualNumber, y::DualNumber) = x * inv(y)
25 end;
```

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

```
1  function Base.log(x::DualNumber{T}) where {T}
2      DualNumber(log(x.value), x.deriv / x.value)
3  end
```

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),
8      (error = err,)
9  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
```

The ForwardDiff.jl package provides a full implementation of these facilities.

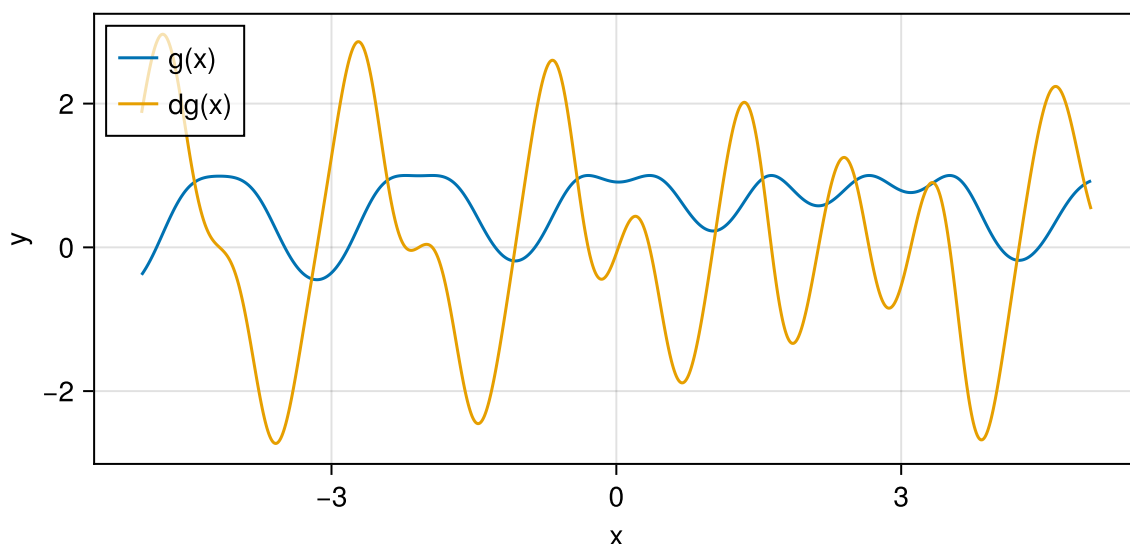testdual1 (generic function with 1 method)

```
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)
```

dg (generic function with 1 method)

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



```
1  let
2      X = (-5:0.01:5)
3      fig = Figure(; size = (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; backgroundcolor = RGBA(1.0, 1.0, 1.0, 0.7), position = :lt)
8      fig
9  end
```

It is also possible to use ForwardDiff with DifferentiationInterface.

(353.0, [42.0, 23.0])

```
1  value_and_gradient(f, AutoForwardDiff(), [x, y])
```

## Forward or Reverse ?

Generally, for the calculation of the Jacobian of a function $f : \mathbb{R}^n \to \mathbb{R}^m$, forward mode AD is more efficient if $n << m$ and reverse mode is more efficient if $n >> m$. The later is important for the calculation of the gradients of the loss functions with respect to the weights in neural network training.

In this course we will stick to the former, as the Jacobian calculation for PDE discretization can be subdivided into operations e.g. on each finite element which have a small number $n$ of input variables.

# 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 nonlinar operator $A : D \to \mathbb{R}^n$ 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 $u$, $M(u) : \mathbb{R}^n \to \mathbb{R}^n$ 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
```

# 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,

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

contraction (generic function with 1 method)

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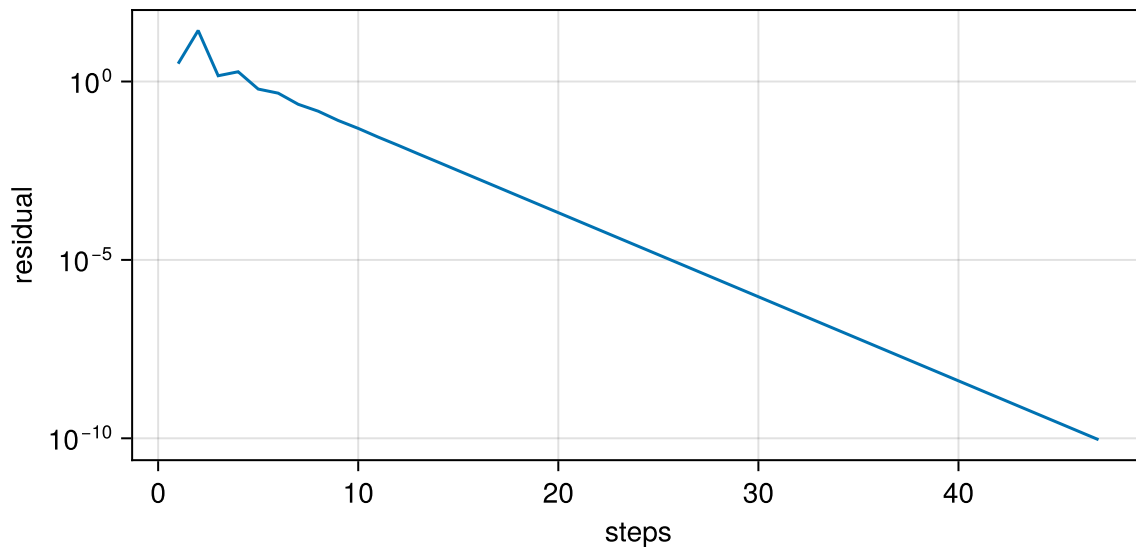
```
 1  function plothistory(history::Vector{<:Number})
 2      fig = Figure(; size = (600, 300))
 3      axis = Axis(fig[1, 1]; xlabel = "steps", ylabel = "residual",
 4                  yscale = log10)
 5      lines!(axis, 1:length(history), history)
 6      #   axislegend(axis,bgcolor = RGBA(1.0,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

[8.50882, 0.0538958, 1.29456, 0.327268, 0.76749, 0.487702, 0.640077, 0.548586, 0.60068

```
 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

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

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 a 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:

- Potenially 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:

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

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

## 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
```

newton (generic function with 1 method)
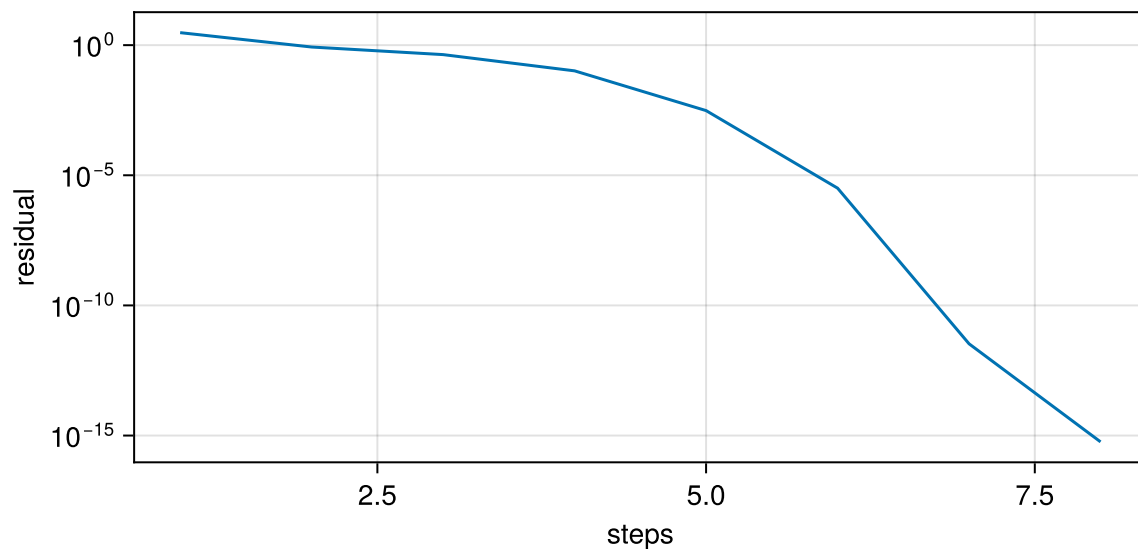
```
1  function newton(A, b, u0; tol = 1.0e-12, maxit = 100)
2      history = Float64[]
3      u = copy(u0)
4      it = 0
5      converged = false
6      while !converged && it < maxit
7          res = A(u) - b
8          jac = ForwardDiff.jacobian((v) -> A(v) - b, u)
9          h = jac \ res
10         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, 3.3

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



```
1  plothistory(newton_history1)
```

Calculate function and derivative at once ?

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

A2λ (generic function with 1 method)

```
1  function A2λ(x, λ)
2      [x[1] + 10λ * x[1]^5 + 3λ * x[2] * x[3],
3        0.1 * x[2] + 10λ * x[2]^5 - 3λ * x[1] - x[3],
4        10λ * x[3]^5 + 10λ * x[1] * x[2] * x[3] + x[3] / 100]
5  end
```

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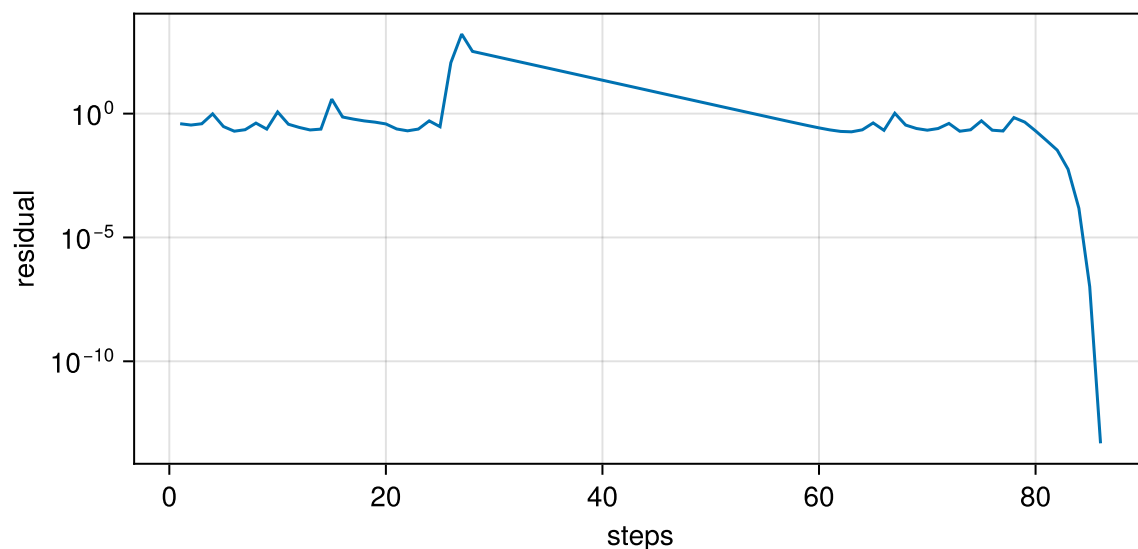
```
1  U02 = [1, 1.0, 1.0]
```

([-0.188484, 0.198519, 0.488388], [0.39077, 0.345694, 0.389908, 0.977557, 0.300465, 0.1

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

[0.0, 8.32667e-17, -5.55112e-17]

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

Newton steps: 86



```
1  plothistory(hist2)
```

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

## Damped Newton scheme

There are may 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,
 2                   damp_growth = 1)
 3      result = DiffResults.JacobianResult(u0)
 4      history = Float64[]
 5      u = copy(u0)
 6      it = 1
 7      while it < maxit
 8          ForwardDiff.jacobian!(result, (v) -> A(v) - b, u)
 9          res = DiffResults.value(result)
10          jac = DiffResults.jacobian(result)
11          h = jac \ res
12          u .-= damp * h
13          nm = norm(h)
14          push!(history, nm)
15          if nm < tol
16              return u, history
17          end
18
19          it = it + 1
20          damp = min(damp * damp_growth, 1.0)
21      end
22      throw("convergence failed")
23  end
```
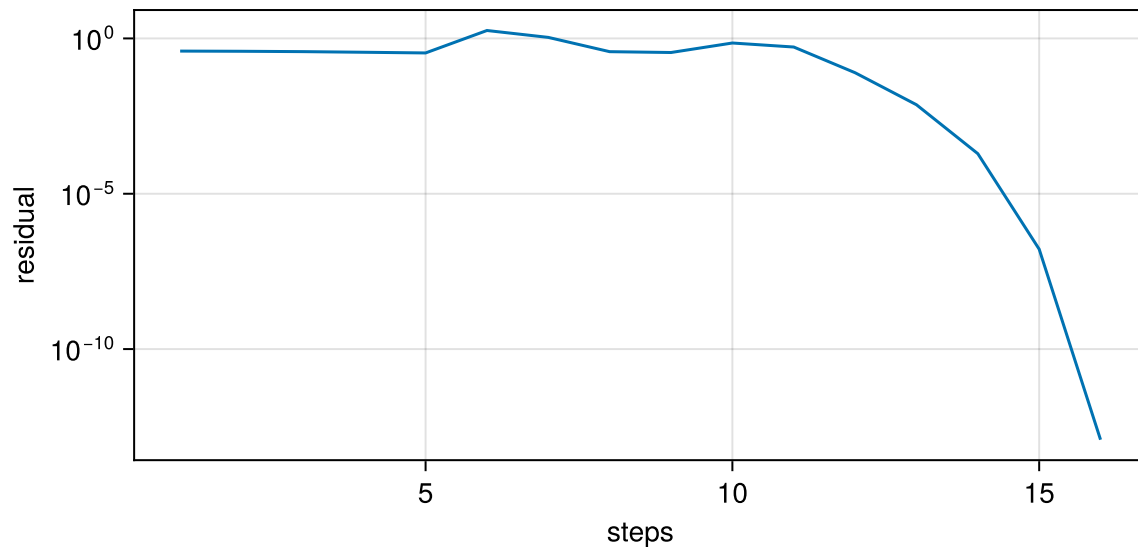
In this implementation, we also try to save work by evaluating result and Jacobian once.

```
([-0.188484, 0.198519, 0.488388], [0.39077, 0.38541, 0.375394, 0.358292, 0.340649, 1.79
```

```
 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 convergece behaviour. If we would keep the damping parameter less than 1, we loose 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

embed_newton (generic function with 1 method)

```
 1  function embed_newton(A, F, U0; δ0 = 0.1, δgrowth = 1.2, λ0 = 0, λ1 = 1)
 2      U = copy(U0)
 3      allhist = Vector[]
 4      λ = λ0
 5      δ = δ0
 6      while true
 7          U, hist = newton(x -> A(x, λ), F, U)
 8          push!(allhist, hist)
 9          if λ == λ1
10              break
11          end
12          λ = min(λ + δ, λ1)
13          δ *= δgrowth
14      end
15      U, allhist
16  end
```

([-0.188484, 0.198519, 0.488388], [[100.408, 1.41554e-14], [28.0258, 16.6762, 13.3379,

```
 1  res4, hist4 = embed_newton(A2λ, F2, U02; δ0 = 0.01, δgrowth = 5.0)
```

[0.0, 8.32667e-17, -5.55112e-17]

```
 1  A2λ(res4, 1.0) - F2
```
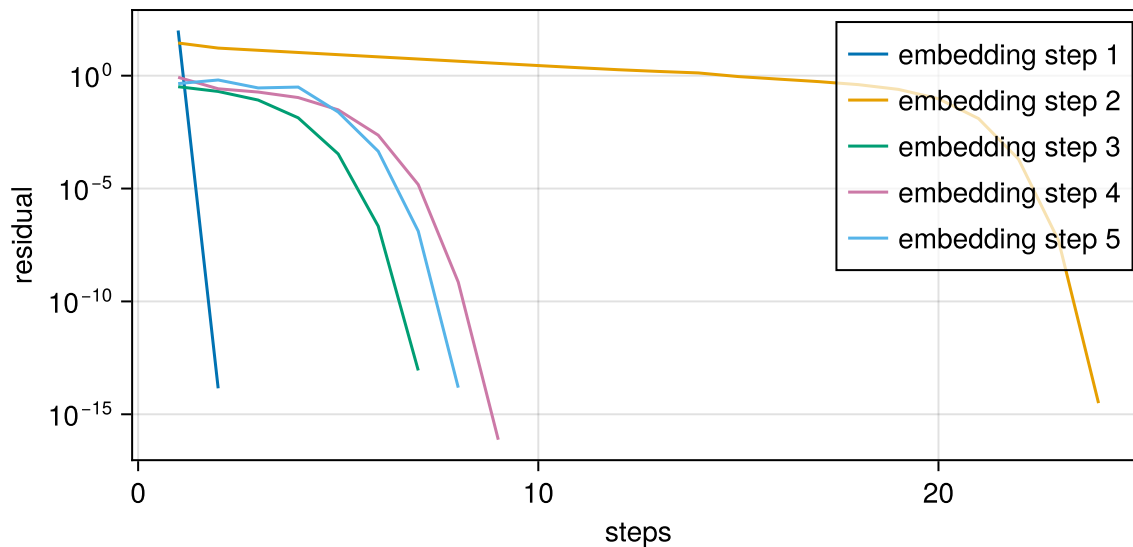
Newton steps: 50

plothistory (generic function with 2 methods)

```
 1  function plothistory(allhistory::Vector{<:AbstractVector})
 2      fig = Figure(; size = (600, 300))
 3      axis = Axis(fig[1, 1]; xlabel = "steps", ylabel = "residual",
 4                  yscale = log10)
 5      i = 1
 6      for history in allhistory
 7          lines!(1:length(history), history; label = "embedding step $(i)")
 8          i = i + 1
 9      end
10      axislegend(axis; backgroundcolor = RGBA(1.0, 1.0, 1.0, 0.7), position = :rt)
11      fig
12  end
```

```
1 plothistory(hist4)
```

# NLsolve.jl

```
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: [0.0575824474961491, 0.48399543065156675, 0.04126490344485912]
          * 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): 74
          * Jacobian Calls (df/dx): 37
```

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

```
[0.0175049, -2.60121e-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: [-0.18848435787886608, 0.19851914493466086, 0.48838826112819433]
          * Inf-norm of residuals: 0.000000
          * Iterations: 87
          * Convergence: true
            * |x - x'| < 0.0e+00: false
            * |f(x)| < 1.0e-08: true
          * Function Calls (f): 88
          * Jacobian Calls (df/dx): 88
```

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

```
[-5.3163e-12, 3.84554e-13, 6.42428e-11]
```

```
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: [-0.18848435786937284, 0.19851914494226675, 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,
2                   autodiff = :forward)
```

```
[-7.93809e-15, 4.16334e-16, 1.06332e-13]
```

```
1  A2λ(nlres3.zero, 1.0) - F2
```

# Summary

- Automatic differentiation significantly simplifies handling of nonlinear systems
- 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
```