

# Automatic differentiation

Note Title

2025-03-28

We have code to compute  $y = f(x)$   $f: \mathbb{R}^m \rightarrow \mathbb{R}^n$

Can we also compute  $f'(x)$  without having to write code for it (and obtain formulas).

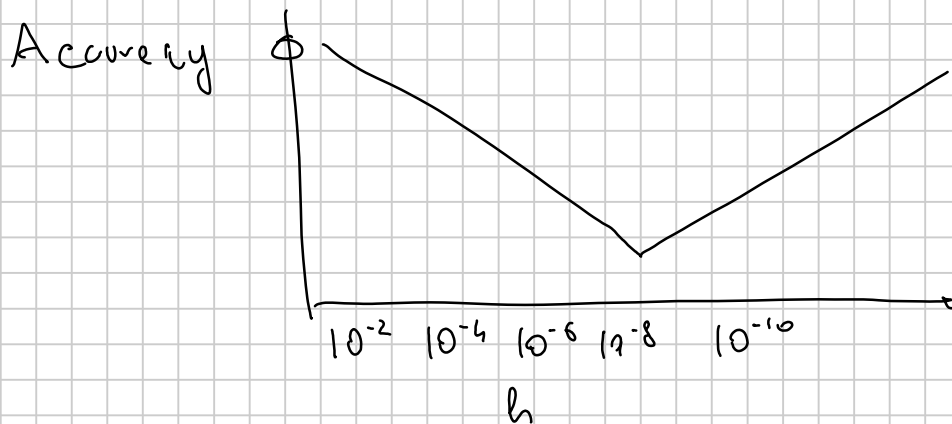
---

First classical approach: numerical differentiation

• set  $h > 0$

$f: \mathbb{R} \rightarrow \mathbb{R}$

• compute  $g = \frac{f(x+h) - f(x)}{h}$



Analytical error:

$$f(x+h) = f(x) + f'(x)h + \frac{1}{2}f''(\xi)h^2$$

$$\frac{f(x+h) - f(x)}{h} - f'(x) = \frac{1}{2}f''(\xi)h$$

$$|g - f'(x)| = \frac{1}{2}f''(\xi)h$$

Machine error:

$$f(x)(1+\epsilon_1)$$

$$|\epsilon_1| \leq u$$

even with the best

$$f(x+h)(1+\epsilon_2)$$

$$|\epsilon_2| \leq u$$

algorithm in the world

$$\tilde{g} = \frac{f(x+h)(1+\epsilon_1) - f(x)(1+\epsilon_2)}{h}$$

$$|\tilde{g} - g| \leq \left| \frac{1}{h} f(x+h) \right| u + \left| \frac{1}{h} f(x) \right| u$$

$$|\tilde{g} - f'(x)| \leq |\tilde{g} - g| + |g - f'(x)| \leq \frac{1}{2} f''(\xi) h + \frac{|f(x+h)| + |f(x)|}{h} u$$

minimum error when

$$\frac{1}{2} f''(\xi) h = \frac{|f(x+h)| + |f(x)|}{h} u$$

$h \sim u^{\frac{1}{2}}$  gives error  $u^{\frac{1}{2}}$

Ex: do the same analysis for

$$g = \frac{f(x+h) - f(x-h)}{2h}$$

$$\text{error} \approx u^{2/3} \sim 10^{-10} - 10^{-11}$$

Jobs: we can get better, if we write our function to support also matrix arguments.

then,  $f(J_A) = \begin{bmatrix} f(A) & f'(A) & \dots & \frac{1}{(k-1)!} f^{(k-1)}(A) \\ \circ & & & \\ & & & \\ & & & f(A) \end{bmatrix}$

### Automatic differentiation

- full machine precision

- works for a wide variety of constructs

Change code so that it works on matrix arguments:

```
z=x
while |z| < 5
    z = z + 2
end
```

=>

```
z=x
while |z(1,1)| < 5
    z = z + 2 * eye(n)
end
```

$$y = e^z$$

$$y = \text{expm}(z)$$

If we want to compute the function on a Jordan block of size  $n=3$ , all the matrices we encounter are of the form

$$\begin{bmatrix} a & b & c \\ 0 & a & b \\ 0 & 0 & a \end{bmatrix} = a \cdot I + b \cdot E + c \cdot E^2 \quad E = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix}$$

We can treat  $E$  as an indeterminate, and compute

$$f(s+E) \quad \text{in the polynomial field} \quad \frac{\mathbb{R}[E]}{(E^3)}$$

This lowers the cost from  $O(n^3)$   
to  $O(n^2)$  (naive polynomial arithmetic)  
 $O(n \log n)$  (FFT)

From an analysis point of view:

start computing from  $x = 5 + \epsilon \quad \leftarrow [5 \ 1 \ 0]$

$$w = x^2 = (5 + \epsilon)(5 + \epsilon) = 25 + 10\epsilon + \epsilon^2 + O(\epsilon^3) \quad \leftarrow [5 \ 1 \ 0] * [5 \ 1 \ 0] = [25 \ 10 \ 1]$$

$$z = x + 5 = 10 + \epsilon + O(\epsilon^3)$$

$$y = wz = (25 + 10\epsilon + \epsilon^2)(10 + \epsilon) = 250 + 125\epsilon + 20\epsilon^2 + O(\epsilon^3)$$

The result gives the power series expansion of  $y$  in terms of  $\epsilon$ , and this reveals derivatives.

---

We can implement this in code using object-oriented programming.

Idea: define new types of variables (classes)  
that contain data (members)

and functions/operations that operate on them (methods)

A Taylor expansion of degree 2 can be represented by

$$a \text{ vector: } 5 + \epsilon \quad \rightarrow \quad [5 \ 1 \ 0]$$

Operations to define:

- $[a_0 \ a_1 \ a_2] * [b_0 \ b_1 \ b_2] = [a_0 b_0 \ a_1 b_0 + a_0 b_1 \ a_2 b_0 + a_1 b_1 + a_0 b_2]$
- $[a_0 \ a_1 \ a_2] + [b_0 \ b_1 \ b_2] = [a_0 + b_0 \ a_1 + b_1 \ a_2 + b_2]$
- $[a_0 \ a_1 \ a_2] + K = [a_0 \ a_1 \ a_2] + [K \ 0 \ 0]$  for  $K \in \mathbb{R}$ .

Once these are defined,  $f(\text{Taylor}([5 \ 1 \ 0]))$  returns  $f$  with its derivatives.

$\exp(\text{Taylor})$  can be defined via

$$\exp(a + b\varepsilon + c\varepsilon^2) = \sum_{k \geq 0} \frac{1}{k!} (a + b\varepsilon + c\varepsilon^2)^k$$

If there is only one derivative:

$a$  with derivative  $a'$

$$b = \exp(a) \quad b' = \exp(a) a'$$

Generally:

$$z = f(a, b, \dots)$$

$$z' = \frac{\partial f}{\partial a} a' + \frac{\partial f}{\partial b} b' + \dots$$

$$z'' = \frac{\partial^2 f}{\partial a^2} (a')^2 + \frac{\partial^2 f}{\partial a \partial b} a' b' + \frac{\partial^2 f}{\partial b^2} (b')^2 + \frac{\partial f}{\partial b} b'' + \dots$$

More complicated expressions for further derivatives.

The  $n=2$  (only one derivative) case is known as the dual numbers:

The dual numbers are the ring  $\frac{\mathbb{R}[\varepsilon]}{(\varepsilon^2)}$

$$\text{e.g. } (5+\epsilon)(2+3\epsilon) = 10 + 17\epsilon$$

With multivariate functions, this still works:

$$z = f(a, b, \dots)$$

$$z' = \frac{\partial f}{\partial a} a' + \frac{\partial f}{\partial b} b'$$

matrix multiplication  
Jacobians

$\frac{dx}{dx} = I$ , then for each intermediate variable  $z$  compute  $\frac{dz}{dx} \in \mathbb{R}^{n_z \times n_x}$

$$c = a \cdot b \quad \frac{dc}{dx} = a \frac{db}{dx} + \frac{da}{dx} \cdot b$$

The cost depends on the dimensions of the involved matrix multiplications

If  $x \in \mathbb{R}^n$ ,  $m$  intermediate variables are of dimension  $n$ ,  
 $y \in \mathbb{R}^m$ ,  $m < n$ , the cost is  $O(n^3)$

This method is more effective when  $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ ,  
with  $n \leq m$ , and all intermediate variables of dimension  $n$ .

Neural networks: a big fitting problem  $\min L(w)$

where  $w \in \mathbb{R}^n$ ,  $n$  very large,  $L: \mathbb{R}^n \rightarrow \mathbb{R}$ .

Since  $n$  (and intermediate variables) are very large, this is inefficient!

Idea to circumvent the problem: start from the last variable: "reverse mode" of automatic differentiation.

Take our model problem

$$w = x^2 \quad \text{input: } x$$

$$z = x + 5 \quad \text{output: } y$$

$$y = w \cdot z$$

In the usual "forward mode", we compute derivatives w.r.t.  $x$  for each variable:

$$\frac{dx}{dx} = 1, \quad \frac{dw}{dx}, \quad \frac{dz}{dx}, \quad \frac{dy}{dx}$$

In reverse mode, we try to start from the end:

$$\frac{dy}{dy} = 1, \quad \frac{dy}{dw}, \quad \frac{dy}{dz}, \quad \frac{dy}{dx}$$

(In a multivariate context where  $x \in \mathbb{R}^n$ ,  $y \in \mathbb{R}^m$ ,  $n \gg m$ , these are "small"  $m \times \text{something}$  matrices)

