

Differenziazione automatica

Note Title

2021-03-16

Opzioni: 1) derivate simboliche: uso calcolo simbolico/formule

2) Derivate numeriche: scelto $h > 0$, approssimo $f'(x)$ con $f'(x) \approx \frac{f(x+h) - f(x)}{h}$

Due fonti di errore

a) f non è esattamente $f'(x)$:

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

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

b) aritmetica di macchine in $f(x+h)$, $f(x)$

Anche con arit. perfetto

$|\delta_i| \leq u$ prec. macchina

$$\begin{aligned} & \frac{f(x+h)(1+\delta_1) - f(x)(1+\delta_2)}{h} (1+\delta_3)(1+\delta_4) - f \\ &= \frac{f(x+h)(\delta_1 + \delta_2 + \delta_3) + f(x)(\delta_2 + \delta_3 + \delta_4) + O(u^2)}{h} \\ &= \left(\left| \frac{f(x+h)}{h} \right| + \left| \frac{f(x)}{h} \right| \right) u \end{aligned}$$

Errore globale,

$$|\tilde{f}' - f'(x)| \leq \underbrace{\left| \frac{1}{2} f''(\xi) \right|}_{=O(h)} \cdot h + \underbrace{\left(\left| \frac{f(x+h)}{h} \right| + \left| \frac{f(x)}{h} \right| \right)}_{=O(1)} \cdot \frac{u}{h}$$

$u \approx 10^{-16}$

$$h = 10^{-4}$$

$$10^{-4}$$

$$10^{-12}$$

$$h = 10^{-8}$$

$$10^{-8}$$

$$10^{-8}$$

← $h = O(u^{1/2})$

$$h = 10^{-12}$$

$$10^{-12}$$

$$10^{-4}$$

Più formalmente,

$$\frac{\left| \frac{1}{2} f''(\xi) \right| h + \left(\left| \frac{f(x+h)}{h} \right| + \left| \frac{f(x)}{h} \right| \right) \cdot \frac{u}{h}}{2} \geq \sqrt{\left| \frac{1}{2} f''(\xi) \right| \left(\left| \frac{f(x+h)}{h} \right| + \left| \frac{f(x)}{h} \right| \right) \cdot \frac{u}{h}}$$

$$= C \cdot u^{\frac{1}{2}}, \text{ con uguaglianza se } \left| \frac{1}{2} f''(\xi) \right| h = \left(|f(x+h)| + |f(x)| \right) \frac{u}{h}$$

Complex-step: supporto f domofo, e di avere codice per calcolare anche con input complessi

uso incremento ih immag. puro

$$f(x+ih) = f(x) + f'(x) \cdot ih - \frac{1}{2} f''(x) h^2 + O(h^3)$$

$$\operatorname{Im} \left[\frac{f(x+ih) - f(x)}{h} \right] - f'(x) = \frac{1}{2} f''(x) h + O(h^2)$$

$$|\tilde{a} - f'(x)| \leq \underbrace{\left| \frac{1}{6} f'''(\xi) \right| h^2}_{O(h^2)} + \underbrace{\left| \operatorname{Im} f(x+ih) \right| \frac{u}{h}}_{O(h)} \leftarrow O(h^2) + O(u)$$

$$f(x+ih) = \operatorname{Re} f(x+ih) + \operatorname{Im} f(x+ih) (1+\delta i)$$

$$\operatorname{Im} f(x+ih) (1+\delta i)$$

tipicamente piccola, dell'ordine di h

e anche calcolata (solitamente) con precisione dell'ordine di h esplicitamente,

$$(x+ih)(x+ih) = (x^2 - h^2) + i(x \cdot h + h \cdot x)$$

codice che funziona su hpi più generali

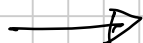
⇒ errore minore.

Se f funziona anche su matrici, cose che richiedono piccoli cambiamenti, ad es.

$$z = x * x$$

$$w = x + 5$$

$$y = z * w$$



$$n = \text{size}(x, 1)$$

$$z = x * x$$

$$w = x + 5 * \text{eye}(n)$$

$$y = z * w$$

Funzione anche con loop, funzioni definite su metodi come l'esponenziale ($\exp(x)$), ...

Questo, in pratica, è propagare sviluppi di Taylor:

invece di

$$\begin{array}{l} z = x * x \\ w = x + 5 \\ y = z * w \end{array} \quad \left| \quad \begin{array}{l} x = 5 \\ z = 25 \\ w = 10 \\ y = 250 \end{array} \right.$$

calcolo w sviluppi di Taylor di una perturb.

$$x = 5 + 1 \cdot \epsilon = 5 + 1 \cdot \epsilon + 0 \cdot \epsilon^2 + O(\epsilon^3)$$

$$z = (5 + 1 \cdot \epsilon + 0 \cdot \epsilon^2 + O(\epsilon^3))(5 + 1 \cdot \epsilon + 0 \cdot \epsilon^2 + O(\epsilon^3)) = 25 + 10 \cdot \epsilon + 1 \cdot \epsilon^2 + O(\epsilon^3)$$

$$w = 5 + 1 \cdot \epsilon + 0 \cdot \epsilon^2 + O(\epsilon^3) + 5 = 10 + 1 \cdot \epsilon + 0 \cdot \epsilon^2 + O(\epsilon^3)$$

$$y = (25 + 10 \epsilon + 1 \epsilon^2 + O(\epsilon^3)) * (10 + 1 \cdot \epsilon + 0 \cdot \epsilon^2 + O(\epsilon^3))$$

$$= 250 + 125 \epsilon + 20 \epsilon^2 + O(\epsilon^3)$$

$$\begin{array}{ccc} \uparrow & \uparrow & \uparrow \\ f(5) & f'(5) & f''(5) \end{array}$$

$$\begin{pmatrix} 25 & 10 & 1 \\ 0 & 25 & 0 \\ 0 & 0 & 25 \end{pmatrix} \cdot \begin{pmatrix} 10 & 1 & 0 \\ 0 & 10 & 1 \\ 0 & 0 & 10 \end{pmatrix}$$

$$x = \text{Taylor}[5 \ 1 \ 0]$$

$$z = x * x = \text{Taylor}[5 \ 1 \ 0] * \text{Taylor}[5 \ 1 \ 0] = \text{Taylor}[25 \ 10 \ 1]$$

$$w = x + 5 = \text{Taylor}[5 \ 1 \ 0] + 5 = \text{Taylor}[10 \ 0 \ 0]$$

$$y = \text{Taylor}[\dots] * \text{Taylor}[\dots]$$

Regola: $\text{Taylor}[a_1, a_2, a_3] * \text{Taylor}[b_1, b_2, b_3]$

$$= \text{Taylor}[a_1 b_1, a_1 b_2 + a_2 b_1, a_1 b_3 + a_2 b_2 + a_3 b_1]$$

Regola: $\text{Taylor}[a_1, a_2, a_3] + \text{Taylor}[b_1, b_2, b_3] = \text{Taylor}[a_1 + b_1, a_2 + b_2, a_3 + b_3]$

Regola: a si può convertire in $\text{Taylor}[a, 0, 0]$

Nello stesso modo, uno può definire Taylor[a]/Taylor[b]

$\exp(\text{Taylor}[a]), \dots$

$$\exp(a_0 + \epsilon a_1 + \epsilon^2 a_2) \approx 1 + (a_0 + \epsilon a_1 + \epsilon^2 a_2) + \frac{1}{2} (a_0 + \epsilon a_1 + \epsilon^2 a_2)^2 + \dots$$

In generale, per una certa op. elementare $z = \varphi(a, b, \dots)$

$$z' = \frac{\partial \varphi}{\partial a}(a, b) a' + \frac{\partial \varphi}{\partial b}(a, b) b'$$

$$z'' = \left(\frac{\partial^2 \varphi}{\partial a^2}(a, b) a' + \frac{\partial^2 \varphi}{\partial a \partial b}(a, b) b' \right) a' + \frac{\partial \varphi}{\partial a}(a, b) a'' + \text{(der. seconde derivate)}$$

Caso più frequente: 1 derivata sola, $[a, a_2] \approx a_1 + \epsilon a_2$

"dual numbers": potete pensare a queste operazioni come definite su un anello $\mathbb{R}[\epsilon]/(\epsilon^2)$

$$(a_1 + a_2 \epsilon)(b_1 + b_2 \epsilon) = a_1 b_1 + (a_1 b_2 + a_2 b_1) \epsilon + \cancel{a_2 b_2 \epsilon^2}$$

(Solo un altro formalismo per le stesse operazioni)

La tecnica usata finora è nota come "forward mode" della "automatic differentiation"

C'è anche un "reverse mode", "back-propagation"

Idea: se il forward mode (per una derivata sola)

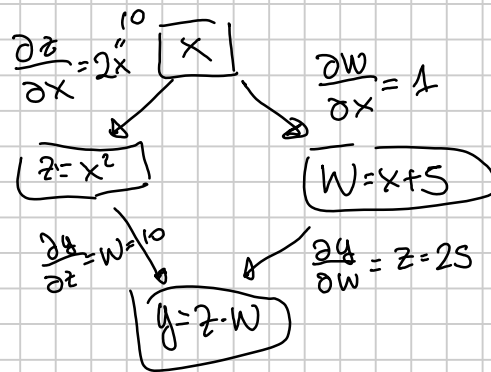
parte da $\frac{dx}{dx} = 1$, e calcola $\frac{da}{dx}$ per ogni quantità calcolata

che compare nel codice $\frac{dz}{dx}$, $\frac{dw}{dx}$, $\frac{dy}{dx}$...

il reverse mode parte dall'output y , con $\frac{\partial y}{\partial y} = 1$,

e calcola $\frac{\partial y}{\partial a}$ per ogni quantità calcolata a che compare nel codice, o partire dall'ultima,

Di solito, tramite grafi: $x=5$



Posso calcolare (dal fondo) derivate parziali di y rispetto a qualunque cosa:

$$\frac{\partial y}{\partial z} = w = 10 \quad \frac{\partial y}{\partial w} = z = 25$$

$$\frac{\partial y}{\partial x} = \frac{\partial y}{\partial w} \cdot \frac{\partial w}{\partial x} + \frac{\partial y}{\partial z} \cdot \frac{\partial z}{\partial x} = 25 \cdot 1 + 10 \cdot 10 = 125$$

Si può costruire algebricamente via trasformazioni del codice, ma è più complicato del forward mode per cui basta definire un tipo.

Qual è il vantaggio del reverse mode? Per $f: \mathbb{R} \rightarrow \mathbb{R}$, nessuno.

Per $f: \mathbb{R}^m \rightarrow \mathbb{R}^n$ di cui voglio tutte le derivate (Jacobiano),

il 'forward mode' richiede meno operazioni se $n \gg m$

il 'reverse mode' richiede meno operazioni se $m \gg n$

Difetti, se $f: \mathbb{R} \rightarrow \mathbb{R}^n$ (nel forward mode) m : basta usare "sviluppi di Taylor per vettori,"

$$\text{cioè, } \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} \rightarrow \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} + \epsilon \begin{bmatrix} 4 \\ 5 \\ 6 \end{bmatrix} = \text{Taylor} \left(\begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}, \begin{bmatrix} 4 \\ 5 \\ 6 \end{bmatrix} \right)$$

Se ho tanti input, devo ripetere il calcolo tante volte considerando diverse perturbazioni: $X = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix}$, devo rispetto a x_1 prima,

$$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix} + \epsilon \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \text{ poi } x_2, \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix} + \epsilon \begin{bmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix}, \dots \begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix} + \epsilon \begin{bmatrix} 0 \\ \vdots \\ 1 \end{bmatrix}$$

Machine learning (ad alto livello): fittinge funzioni con multi input con tecniche di ottimizzazione tipo gradient descent, $X_{k+1} \leftarrow X_k + \frac{\partial f}{\partial X} \cdot \delta_k$