

Automatic differentiation

Mauricio A. Álvarez

Machine Learning and Adaptive Intelligence
The University of Sheffield



The
University
Of
Sheffield.

Contents

Derivatives and ways to compute them

AD modes

- Forward mode
- Reverse mode

Implementations

Derivatives

- Derivatives are required to perform optimisation in several ML algorithms.
- For example, computing the gradient is necessary for **batch gradient descent** and **SGD**.
- Derivatives are also necessary for computing **Hessians** which are used in second-order optimisation methods.

Methods to compute derivatives in computer programs

- Manually working out derivatives and coding them.
- Numeric differentiation using finite difference approximations.
- Symbolic differentiation.
- Automatic differentiation (or algorithmic differentiation).

Example

- Suppose we have the following function

$$f(x, y) = x^2y + y + 2.$$

- We require to compute the gradient of this function, for example, because we want to use it in gradient descent,

$$\frac{df(x, y)}{dz} = \begin{bmatrix} \frac{\partial f(x, y)}{\partial x} \\ \frac{\partial f(x, y)}{\partial y} \end{bmatrix},$$

where $\mathbf{z} = [x \ y]^\top$. \therefore $\frac{df(x, y)}{\partial z}$: derivatives for $f(x, y)$ wrt \mathbf{z} :

$$\mathbf{z} = \begin{bmatrix} x \\ y \end{bmatrix}$$

Manual differentiation (I)

- ❑ We use our calculus knowledge to derive the proper equation.
↙
- ❑ For the function we saw before, we need to apply the following rules of calculus
 - The derivative of a constant is 0.
 - The derivative of ax with respect to x is a , where a is a constant.
 - The derivative of x^a is ax^{a-1} .
 - The derivative is a linear operation so, the derivative of the sum of two functions is the sum of the derivatives.
 - The derivative of a constant times a function, is equal to the constant times the derivative of that function.
- ❑ Using these rules we get the following partial derivatives,

Manual differentiation (II)

- Partial derivative of $f(x, y)$ with respect to x

$$\frac{\partial}{\partial x} f(x, y) = \frac{\partial}{\partial x} (x^2 y + y + 2) = 2xy. \quad \checkmark$$

- Partial derivative of $f(x, y)$ with respect to y

$$\frac{\partial}{\partial y} f(x, y) = \frac{\partial}{\partial y} (x^2 y + y + 2) = x^2 + 1. \quad \checkmark$$

- We can then write

$$\frac{df(x, y)}{dz} = \begin{bmatrix} \frac{\partial f(x, y)}{\partial x} \\ \frac{\partial f(x, y)}{\partial y} \end{bmatrix} = \begin{bmatrix} 2xy \\ x^2 + 1 \end{bmatrix} \quad \checkmark$$

Problems with manual differentiation

Tedious.
易犯错

When a function $f(\cdot)$ depends on many variables or is a rather complicated expression, manual differentiation is **tedious** and prone to mistakes.



Finite difference approximations (I)

- Remember the definition of a derivative of a function $h(x)$ at a point x_0 ,

$$\begin{aligned}\frac{dh(x_0)}{dx} &= \lim_{x \rightarrow x_0} \frac{h(x) - h(x_0)}{x - x_0} \\ &= \lim_{\epsilon \rightarrow 0} \frac{h(x_0 + \epsilon) - h(x_0)}{\epsilon}\end{aligned}$$

- The partial derivative of $h(x, y)$ at point (x_0, y_0) is defined as

x is the derivative

$$\frac{\partial h(x_0, y_0)}{\partial x} = \lim_{\epsilon \rightarrow 0} \frac{h(x_0 + \epsilon, y_0) - h(x_0, y_0)}{\epsilon}$$

y is the derivative

$$\frac{\partial h(x_0, y_0)}{\partial y} = \lim_{\epsilon \rightarrow 0} \frac{h(x_0, y_0 + \epsilon) - h(x_0, y_0)}{\epsilon}$$

Finite difference approximations (II)

easy to implement

```
def f(x, y):
    return x**2*y + y + 2
x_0 = 3
y_0 = 2
epsilon = 1e-6
dfdx_numerical = (f(x_0+epsilon, y_0) - f(x_0, y_0))/epsilon
dfdy_numerical = (f(x_0, y_0+epsilon) - f(x_0, y_0))/epsilon
dfdx_analytical = 2*x_0*y_0
dfdy_analytical = x_0**2 + 1
```

$$f(x, y) = x^2y + y + 2$$
$$\frac{\partial f}{\partial x} = 2xy; \quad \frac{\partial f}{\partial y} = x^2 + 1$$
$$\frac{f(3 + 10^{-6}, y) - f(3, y)}{10^{-6}}$$
$$f(3, y + 10^{-6}) - f(3, y)$$
$$(10^{-6})$$

Script in python for the finite differences

```
In [22]: dfdx_numerical  
Out[22]: 12.000002001855137
```

```
In [23]: dfdx_analytical  
Out[23]: 12
```

$$\frac{\partial f(x, y)}{\partial x}$$

```
In [24]: dfdy_numerical  
Out[24]: 10.000000003174137
```

```
In [25]: dfdy_analytical  
Out[25]: 10
```

$$\frac{\partial f(x, y)}{\partial y}$$

Problems with finite difference approximation

- The result is imprecise and gets worse with more complicated functions.
- We need to call the function at least twice. For big parametric models, we'd need to call the function several times becoming very inefficient.
- The method is easy to implement, so one can use it to test whether the manual implementation is correct.

Symbolic differentiation (I)

- ❑ Symbolic differentiation performs an automatic manipulation of expressions to obtain the corresponding derivative expressions.
- ❑ The mathematical expression is represented using **data structures** (e.g. trees, lists, etc.).
- ❑ It is then possible to follow a mechanistic process to obtain the derivatives.

Symbolic differentiation (II)



Maxima



Maple

Symbolic differentiation with Mathematica

In[33]:= D[x, x]

Out[33]:= 1

In[34]:= D[4 x (1 - x), x]

Out[34]:= 4 (1 - x) - 4 x

In[35]:= D[16 x (1 - x) ((1 - 2 x)^2), x]

Out[35]:= 16 (1 - 2 x)^2 (1 - x) - 16 (1 - 2 x)^2 x - 64 (1 - 2 x) (1 - x) x

In[36]:= D[64 x (1 - x) ((1 - 2 x)^2) ((1 - 8 x + 8 x^2)^2), x]

Out[36]:= 128 (1 - 2 x)^2 (1 - x) x (-8 + 16 x) (1 - 8 x + 8 x^2) +
64 (1 - 2 x)^2 (1 - x) (1 - 8 x + 8 x^2)^2 - 64 (1 - 2 x)^2 x (1 - 8 x + 8 x^2)^2 -
256 (1 - 2 x) (1 - x) x (1 - 8 x + 8 x^2)^2

Problems with symbolic differentiation

- Due to the mechanistic approach, there is usually a lot of **redundancy** in the expressions generated.
- If not handled properly, it produces unnecessary long expressions difficult to make sense of and to evaluate.
- Such behavior is known as *expression swell*.

Example of expression swell

n	I_n	$\frac{dI_n}{dx}$
1	x	1
2	$4x(1-x)$	$4(1-x) - 4x$
3	$16x(1-x)(1-2x)^2$	$16(1-2x)^2(1-x) - 16(1-2x)^2x - 64(1-2x)(1-x)x$
4	$64x(1-x)(1-2x)^2(8x^2 - 8x + 1)^2$	$128(1-2x)^2(1-x)x(-8 + 16x)(1 - 8x + 8x^2) + 64(1-2x)^2(1-x)(1-8x + 8x^2)^2 - 64(1-2x)^2x(1-8x + 8x^2)^2 - 256(1-2x)(1-x)x(1-8x + 8x^2)^2$

Logistic map $I_n = 4I_n(1 - I_n)$, $I_1 = x$.

Expression swell
can be simplified

Simplify with Mathematica

```
In[40]:= D[16 x (1 - x) ((1 - 2 x)^2), x]
```

```
Out[40]= 16 (1 - 2 x)^2 (1 - x) - 16 (1 - 2 x)^2 x - 64 (1 - 2 x) (1 - x) x
```

```
In[39]:= Simplify[16 (1 - 2 x)^2 (1 - x) - 16 (1 - 2 x)^2 x - 64 (1 - 2 x) (1 - x) x]
```

```
Out[39]= -16 (-1 + 10 x - 24 x^2 + 16 x^3)
```

```
In[41]:= D[64 x (1 - x) ((1 - 2 x)^2) ((1 - 8 x + 8 x^2)^2), x] Expression swell
```

```
Out[41]= 128 (1 - 2 x)^2 (1 - x) x (-8 + 16 x) (1 - 8 x + 8 x^2) +  
64 (1 - 2 x)^2 (1 - x) (1 - 8 x + 8 x^2)^2 - 64 (1 - 2 x)^2 x (1 - 8 x + 8 x^2)^2 -  
256 (1 - 2 x) (1 - x) x (1 - 8 x + 8 x^2)^2
```

```
In[42]:= Simplify[128 (1 - 2 x)^2 (1 - x) x (-8 + 16 x) (1 - 8 x + 8 x^2) +  
64 (1 - 2 x)^2 (1 - x) (1 - 8 x + 8 x^2)^2 - 64 (1 - 2 x)^2 x (1 - 8 x + 8 x^2)^2 -  
256 (1 - 2 x) (1 - x) x (1 - 8 x + 8 x^2)^2]
```

```
Out[42]= -64 (-1 + 42 x - 504 x^2 + 2640 x^3 - 7040 x^4 + 9984 x^5 - 7168 x^6 + 2048 x^7)
```

Automatic differentiation

- AD is concerned about exact numerical computation of the derivatives, rather than their actual symbolic form.
- It computes the derivative by **only storing the values** of intermediate sub-expressions.
- It uses a combination of: **symbolic differentiation at the elementary operation level** and **keeping intermediate numerical results**.

Contents

Derivatives and ways to compute them

AD modes

- Forward mode
- Reverse mode

Implementations

Evaluation trace

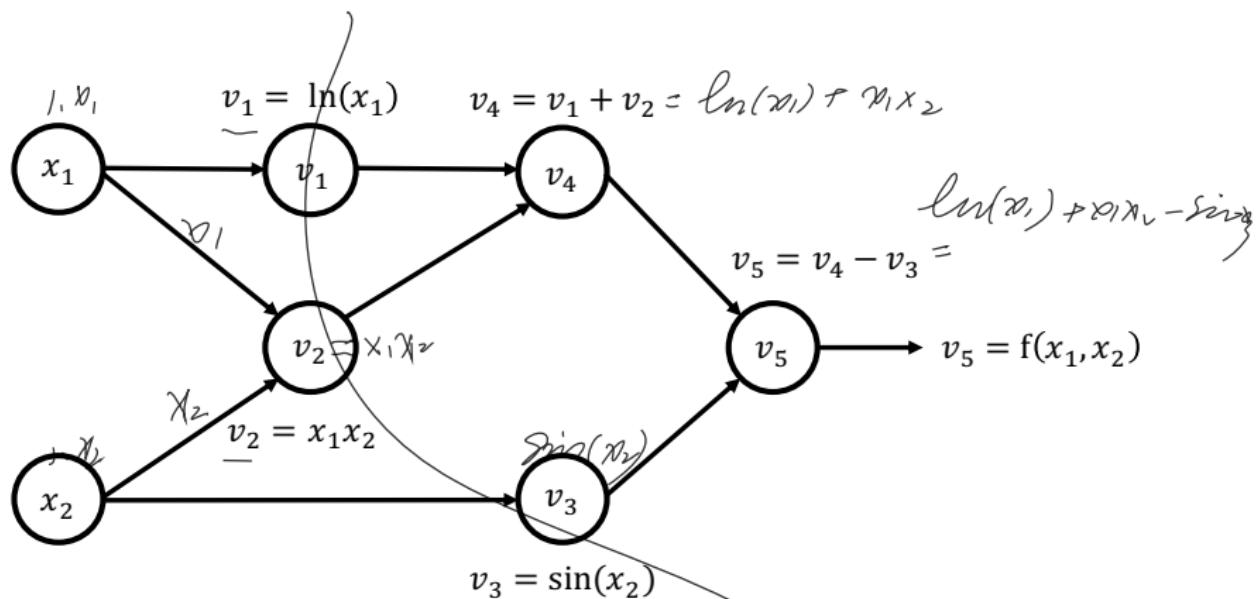
- *Evaluation trace*: composition of elementary operations that lead to a full expression.
- As an example, let us consider the function

$$f(x_1, x_2) = \ln(x_1) + x_1 x_2 - \sin(x_2).$$

- The inputs are x_1 and x_2 . $v_1 = \ln(x_1)$
- The elementary operations include $v_2 = v_1 \cdot v_2$
 - $v_1 = \ln(x_1)$
 - $v_2 = x_1 x_2$
 - $v_3 = \sin(x_2)$
 - $v_4 = v_1 + v_2$
 - $v_5 = v_4 - v_3$
 - $f(x_1, x_2) = v_5$. $v_4 = v_1 + v_2$ $v_8 = v_4 - v_3$ $f(x_1, x_2) = v_5$

Computational graph

$$f(x_1, x_2) = \ln(x_1) + x_1 x_2 - \sin(x_2)$$



General notation

for example:

$$f(x_1, x_2) = \ln(x_1) + x_1 x_2 - \sin(x_2)$$

$$\begin{aligned} x_1 &\in \mathbb{R}^n \\ x_2 &\in \mathbb{R}^n \end{aligned}$$

In this case: $\begin{cases} n=2 \\ i=1, 2 \end{cases}$

- Let $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$.

$$\begin{cases} v_{1-2} = v_1 = x_1 \\ v_{2-2} = v_2 = x_2 \end{cases} \text{ are input var}$$

- Variables $v_{i-n} = x_i$, where $i = 1, \dots, n$ are the input variables.

$$\begin{cases} v_1 = x_1 \\ v_2 = x_2 \end{cases}$$

$$v_0 = x_0$$

$v_1 = \ln(x_1) \rightarrow$ is an intermediate variable.

- Variables v_i , with $i = 1, \dots, l$ are the intermediate variables.

$$\begin{cases} v_1 = \ln(x_1) \\ v_2 = x_1 x_2 \end{cases} \quad l=? \quad i=1, 2, 3, \Rightarrow l=3$$

$$v_3 = \sin(x_2)$$

$$v_4 = \sin(x_3)$$

$$v_5 = \sin(x_4)$$

$$v_6 = \sin(x_5)$$

$$v_7 = \sin(x_6)$$

$$v_8 = \sin(x_7)$$

$$v_9 = \sin(x_8)$$

$$v_{10} = \sin(x_9)$$

$$v_{11} = \sin(x_{10})$$

$$v_{12} = \sin(x_{11})$$

$$v_{13} = \sin(x_{12})$$

$$v_{14} = \sin(x_{13})$$

$$v_{15} = \sin(x_{14})$$

$$v_{16} = \sin(x_{15})$$

$$v_{17} = \sin(x_{16})$$

$$v_{18} = \sin(x_{17})$$

$$v_{19} = \sin(x_{18})$$

$$v_{20} = \sin(x_{19})$$

$$v_{21} = \sin(x_{20})$$

$$v_{22} = \sin(x_{21})$$

$$v_{23} = \sin(x_{22})$$

$$v_{24} = \sin(x_{23})$$

$$v_{25} = \sin(x_{24})$$

$$v_{26} = \sin(x_{25})$$

$$v_{27} = \sin(x_{26})$$

$$v_{28} = \sin(x_{27})$$

$$v_{29} = \sin(x_{28})$$

$$v_{30} = \sin(x_{29})$$

$$v_{31} = \sin(x_{30})$$

$$v_{32} = \sin(x_{31})$$

$$v_{33} = \sin(x_{32})$$

$$v_{34} = \sin(x_{33})$$

$$v_{35} = \sin(x_{34})$$

$$v_{36} = \sin(x_{35})$$

$$v_{37} = \sin(x_{36})$$

$$v_{38} = \sin(x_{37})$$

$$v_{39} = \sin(x_{38})$$

$$v_{40} = \sin(x_{39})$$

$$v_{41} = \sin(x_{40})$$

$$v_{42} = \sin(x_{41})$$

$$v_{43} = \sin(x_{42})$$

$$v_{44} = \sin(x_{43})$$

$$v_{45} = \sin(x_{44})$$

$$v_{46} = \sin(x_{45})$$

$$v_{47} = \sin(x_{46})$$

$$v_{48} = \sin(x_{47})$$

$$v_{49} = \sin(x_{48})$$

$$v_{50} = \sin(x_{49})$$

$$v_{51} = \sin(x_{50})$$

$$v_{52} = \sin(x_{51})$$

$$v_{53} = \sin(x_{52})$$

$$v_{54} = \sin(x_{53})$$

$$v_{55} = \sin(x_{54})$$

$$v_{56} = \sin(x_{55})$$

$$v_{57} = \sin(x_{56})$$

$$v_{58} = \sin(x_{57})$$

$$v_{59} = \sin(x_{58})$$

$$v_{60} = \sin(x_{59})$$

$$v_{61} = \sin(x_{60})$$

$$v_{62} = \sin(x_{61})$$

$$v_{63} = \sin(x_{62})$$

$$v_{64} = \sin(x_{63})$$

$$v_{65} = \sin(x_{64})$$

$$v_{66} = \sin(x_{65})$$

$$v_{67} = \sin(x_{66})$$

$$v_{68} = \sin(x_{67})$$

$$v_{69} = \sin(x_{68})$$

$$v_{70} = \sin(x_{69})$$

$$v_{71} = \sin(x_{70})$$

$$v_{72} = \sin(x_{71})$$

$$v_{73} = \sin(x_{72})$$

$$v_{74} = \sin(x_{73})$$

$$v_{75} = \sin(x_{74})$$

$$v_{76} = \sin(x_{75})$$

$$v_{77} = \sin(x_{76})$$

$$v_{78} = \sin(x_{77})$$

$$v_{79} = \sin(x_{78})$$

$$v_{80} = \sin(x_{79})$$

$$v_{81} = \sin(x_{80})$$

$$v_{82} = \sin(x_{81})$$

$$v_{83} = \sin(x_{82})$$

$$v_{84} = \sin(x_{83})$$

$$v_{85} = \sin(x_{84})$$

$$v_{86} = \sin(x_{85})$$

$$v_{87} = \sin(x_{86})$$

$$v_{88} = \sin(x_{87})$$

$$v_{89} = \sin(x_{88})$$

$$v_{90} = \sin(x_{89})$$

$$v_{91} = \sin(x_{90})$$

$$v_{92} = \sin(x_{91})$$

$$v_{93} = \sin(x_{92})$$

$$v_{94} = \sin(x_{93})$$

$$v_{95} = \sin(x_{94})$$

$$v_{96} = \sin(x_{95})$$

$$v_{97} = \sin(x_{96})$$

$$v_{98} = \sin(x_{97})$$

$$v_{99} = \sin(x_{98})$$

$$v_{100} = \sin(x_{99})$$

$$v_{101} = \sin(x_{100})$$

$$v_{102} = \sin(x_{101})$$

$$v_{103} = \sin(x_{102})$$

$$v_{104} = \sin(x_{103})$$

$$v_{105} = \sin(x_{104})$$

$$v_{106} = \sin(x_{105})$$

$$v_{107} = \sin(x_{106})$$

$$v_{108} = \sin(x_{107})$$

$$v_{109} = \sin(x_{108})$$

$$v_{110} = \sin(x_{109})$$

$$v_{111} = \sin(x_{110})$$

$$v_{112} = \sin(x_{111})$$

$$v_{113} = \sin(x_{112})$$

$$v_{114} = \sin(x_{113})$$

$$v_{115} = \sin(x_{114})$$

$$v_{116} = \sin(x_{115})$$

$$v_{117} = \sin(x_{116})$$

$$v_{118} = \sin(x_{117})$$

$$v_{119} = \sin(x_{118})$$

$$v_{120} = \sin(x_{119})$$

$$v_{121} = \sin(x_{120})$$

$$v_{122} = \sin(x_{121})$$

$$v_{123} = \sin(x_{122})$$

$$v_{124} = \sin(x_{123})$$

$$v_{125} = \sin(x_{124})$$

$$v_{126} = \sin(x_{125})$$

$$v_{127} = \sin(x_{126})$$

$$v_{128} = \sin(x_{127})$$

$$v_{129} = \sin(x_{128})$$

$$v_{130} = \sin(x_{129})$$

$$v_{131} = \sin(x_{130})$$

$$v_{132} = \sin(x_{131})$$

$$v_{133} = \sin(x_{132})$$

$$v_{134} = \sin(x_{133})$$

$$v_{135} = \sin(x_{134})$$

$$v_{136} = \sin(x_{135})$$

$$v_{137} = \sin(x_{136})$$

$$v_{138} = \sin(x_{137})$$

$$v_{139} = \sin(x_{138})$$

$$v_{140} = \sin(x_{139})$$

$$v_{141} = \sin(x_{140})$$

$$v_{142} = \sin(x_{141})$$

$$v_{143} = \sin(x_{142})$$

$$v_{144} = \sin(x_{143})$$

$$v_{145} = \sin(x_{144})$$

$$v_{146} = \sin(x_{145})$$

$$v_{147} = \sin(x_{146})$$

$$v_{148} = \sin(x_{147})$$

$$v_{149} = \sin(x_{148})$$

$$v_{150} = \sin(x_{149})$$

$$v_{151} = \sin(x_{150})$$

$$v_{152} = \sin(x_{151})$$

$$v_{153} = \sin(x_{152})$$

$$v_{154} = \sin(x_{153})$$

$$v_{155} = \sin(x_{154})$$

$$v_{156} = \sin(x_{155})$$

$$v_{157} = \sin(x_{156})$$

$$v_{158} = \sin(x_{157})$$

$$v_{159} = \sin(x_{158})$$

$$v_{160} = \sin(x_{159})$$

$$v_{161} = \sin(x_{160})$$

$$v_{162} = \sin(x_{161})$$

$$v_{163} = \sin(x_{162})$$

$$v_{164} = \sin(x_{163})$$

$$v_{165} = \sin(x_{164})$$

$$v_{166} = \sin(x_{165})$$

$$v_{167} = \sin(x_{166})$$

$$v_{168} = \sin(x_{167})$$

$$v_{169} = \sin(x_{168})$$

$$v_{170} = \sin(x_{169})$$

$$v_{171} = \sin(x_{170})$$

$$v_{172} = \sin(x_{171})$$

$$v_{173} = \sin(x_{172})$$

$$v_{174} = \sin(x_{173})$$

$$v_{175} = \sin(x_{174})$$

$$v_{176} = \sin(x_{175})$$

$$v_{177} = \sin(x_{176})$$

$$v_{178} = \sin(x_{177})$$

$$v_{179} = \sin(x_{178})$$

$$v_{180} = \sin(x_{179})$$

$$v_{181} = \sin(x_{180})$$

$$v_{182} = \sin(x_{181})$$

$$v_{183} = \sin(x_{182})$$

$$v_{184} = \sin(x_{183})$$

$$v_{185} = \sin(x_{184})$$

$$v_{186} = \sin(x_{185})$$

$$v_{187} = \sin(x_{186})$$

$$v_{188} = \sin(x_{187})$$

$$v_{189} = \sin(x_{188})$$

$$v_{190} = \sin(x_{189})$$

$$v_{191} = \sin(x_{190})$$

$$v_{192} = \sin(x_{191})$$

$$v_{193} = \sin(x_{192})$$

$$v_{194} = \sin(x_{193})$$

$$v_{195} = \sin(x_{194})$$

$$v_{196} = \sin(x_{195})$$

$$v_{197} = \sin(x_{196})$$

$$v_{198} = \sin(x_{197})$$

$$v_{199} = \sin(x_{198})$$

$$v_{200} = \sin(x_{199})$$

$$v_{201} = \sin(x_{200})$$

$$v_{202} = \sin(x_{201})$$

$$v_{203} = \sin(x_{202})$$

$$v_{204} = \sin(x_{203})$$

$$v_{205} = \sin(x_{204})$$
</div

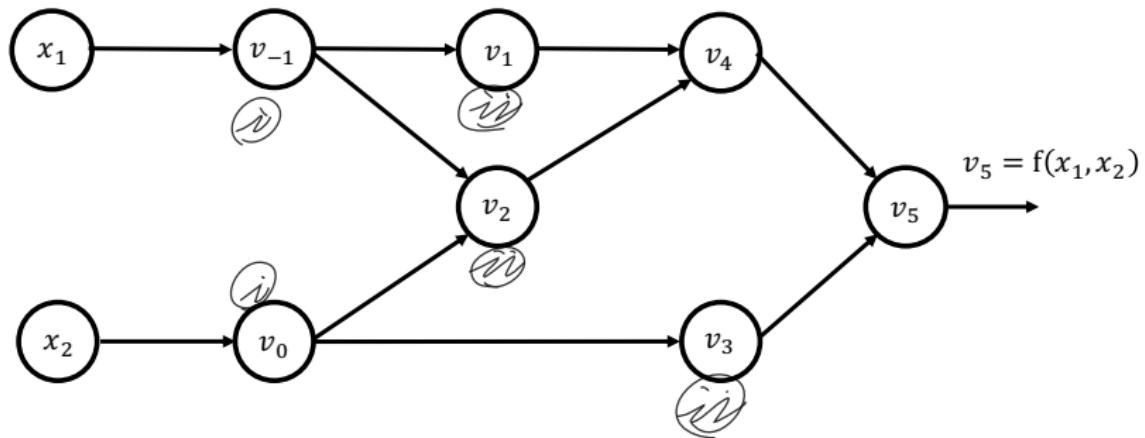
New computational graph

which of these are ...

i) input variables ?

ii) intermediate variables

iii) output variables



Jacobian

a function depends on several variables $f(x_1, x_2, \dots, x_n)$

- Say that we have several functions $y_i = f_i(\cdot)$ for $i = 1, \dots, m$ that depend on several input variables x_1, x_2, \dots, x_n ,

$$y_1 = f_1(x_1, \dots, x_n)$$

$$y_2 = f_2(x_1, \dots, x_n)$$

$$\vdots \quad \vdots$$

$$y_m = f_m(x_1, \dots, x_n)$$

- The Jacobian \mathbf{J} of dimensions $m \times n$ is a matrix with entries $\underline{\mathbf{J}_{ij} = \frac{\partial f_i}{\partial x_j}}$ given as what does the sign mean?

why not use $\left(\frac{\partial f_i}{\partial x_j} \right)$? $\mathbf{J} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \dots & \frac{\partial f_1}{\partial x_n} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & \dots & \frac{\partial f_2}{\partial x_n} \\ \vdots & \vdots & \vdots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \frac{\partial f_m}{\partial x_2} & \dots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}$

$\frac{\partial f_i}{\partial x_j} \rightarrow$ get the derivative of f_i , with respect to x_j .

$$\mathbf{J}_{ij} = \frac{\partial f_i}{\partial x_j}$$

Contents

Derivatives and ways to compute them

AD modes

Forward mode

Reverse mode

Implementations

Forward accumulation mode

- Forward accumulation mode or tangent linear mode.
- To compute the derivative of f with respect to x_1 , each intermediate variable v_i has a derivative, and calculate the derivative of functions (intermediate)
$$\dot{v}_i = \frac{\partial v_i}{\partial x_1}$$
- For each evaluation (or forward primal) trace, it builds a forward derivative (or tangent) trace.
- Essentially, this forward derivative trace is just implementing the chain rule of differentiation

$$\frac{dy}{dx} = \frac{dy}{du} \frac{du}{dz} \frac{dz}{dx}.$$

Forward primal trace and tangent trace for $\frac{\partial y}{\partial x_1}$ (I)

- Let us compute the forward tangent trace $\frac{\partial y}{\partial x_1}$ for the function we had before

$$y = f(x_1, x_2) = \ln(x_1) + x_1 x_2 - \sin(x_2).$$

- The following table shows both the forward primal trace and the forward tangent trace

Forward primal trace	Forward tangent trace
$v_{-1} = x_1$	$\dot{v}_{-1} = \dot{x}_1$
$v_0 = x_2$	$\dot{v}_0 = \dot{x}_2$
$v_1 = \ln v_{-1}$	$\dot{v}_1 = \frac{1}{v_{-1}} \dot{v}_{-1}$
$v_2 = v_{-1} \times v_0$	$\dot{v}_2 = \dot{v}_{-1} \times v_0 + v_0 \times \dot{v}_{-1}$
$v_3 = \sin(v_0)$	$\dot{v}_3 = \dot{v}_0 \times \cos(v_0)$
$v_4 = v_1 + v_2$	$\dot{v}_4 = \dot{v}_1 + \dot{v}_2$
$v_5 = v_4 - v_3$	$\dot{v}_5 = \dot{v}_4 - \dot{v}_3$
$y = v_5$	$\dot{y} = \dot{v}_5$

Forward primal trace and tangent trace for $\frac{\partial y}{\partial x_1}$ (II)

We now compute the derivative $\frac{\partial y}{\partial x_1}$ at $x_1 = 2, x_2 = 5$.

Forward primal trace and tangent trace for $\frac{\partial y}{\partial x_1}$ (II)

We now compute the derivative $\frac{\partial y}{\partial x_1}$ at $x_1 = 2, x_2 = 5$.

Forward primal trace	Forward tangent trace
$v_{-1} = x_1 = 2$	$\dot{v}_{-1} = \dot{x}_1$
$v_0 = x_2 = 5$	$\dot{v}_0 = \dot{x}_2$
$v_1 = \ln v_{-1}$	$\dot{v}_1 = \frac{1}{v_{-1}} \dot{v}_{-1}$
$v_2 = v_{-1} \times v_0$	$\dot{v}_2 = \dot{v}_{-1} \times v_0 + v_0 \times \dot{v}_{-1}$
$v_3 = \sin(v_0)$	$\dot{v}_3 = \dot{v}_0 \times \cos(v_0)$
$v_4 = v_1 + v_2$	$\dot{v}_4 = \dot{v}_1 + \dot{v}_2$
$v_5 = v_4 - v_3$	$\dot{v}_5 = \dot{v}_4 - \dot{v}_3$
$y = v_5$	$\dot{y} = \dot{v}_5$

Forward primal trace and tangent trace for $\frac{\partial y}{\partial x_1}$ (II)

We now compute the derivative $\frac{\partial y}{\partial x_1}$ at $x_1 = 2, x_2 = 5$.

Forward primal trace	Forward tangent trace
$v_{-1} = x_1 = 2$	$\dot{v}_{-1} = \dot{x}_1 = 1$
$v_0 = x_2 = 5$	$\dot{v}_0 = \dot{x}_2 = 0$
$v_1 = \ln v_{-1}$	$\dot{v}_1 = \frac{1}{v_{-1}} \dot{v}_{-1}$
$v_2 = v_{-1} \times v_0$	$\dot{v}_2 = \dot{v}_{-1} \times v_0 + v_0 \times \dot{v}_{-1}$
$v_3 = \sin(v_0)$	$\dot{v}_3 = \dot{v}_0 \times \cos(v_0)$
$v_4 = v_1 + v_2$	$\dot{v}_4 = \dot{v}_1 + \dot{v}_2$
$v_5 = v_4 - v_3$	$\dot{v}_5 = \dot{v}_4 - \dot{v}_3$
$y = v_5$	$\dot{y} = \dot{v}_5$

Forward primal trace and tangent trace for $\frac{\partial y}{\partial x_1}$ (II)

We now compute the derivative $\frac{\partial y}{\partial x_1}$ at $x_1 = 2, x_2 = 5$.

Forward primal trace	Forward tangent trace
$v_{-1} = x_1 = 2$	$\dot{v}_{-1} = \dot{x}_1 = 1$
$v_0 = x_2 = 5$	$\dot{v}_0 = \dot{x}_2 = 0$
$v_1 = \ln v_{-1} = \ln 2$	$\dot{v}_1 = \frac{1}{v_{-1}} \dot{v}_{-1} = \frac{1}{2}(1)$
$v_2 = v_{-1} \times v_0$	$\dot{v}_2 = \dot{v}_{-1} \times v_0 + v_0 \times \dot{v}_{-1}$
$v_3 = \sin(v_0)$	$\dot{v}_3 = \dot{v}_0 \times \cos(v_0)$
$v_4 = v_1 + v_2$	$\dot{v}_4 = \dot{v}_1 + \dot{v}_2$
$v_5 = v_4 - v_3$	$\dot{v}_5 = \dot{v}_4 - \dot{v}_3$
$y = v_5$	$\dot{y} = \dot{v}_5$

Forward primal trace and tangent trace for $\frac{\partial y}{\partial x_1}$ (II)

We now compute the derivative $\frac{\partial y}{\partial x_1}$ at $x_1 = 2, x_2 = 5$.

Forward primal trace	Forward tangent trace
$v_{-1} = x_1 = 2$	$\dot{v}_{-1} = \dot{x}_1 = 1$
$v_0 = x_2 = 5$	$\dot{v}_0 = \dot{x}_2 = 0$
$v_1 = \ln v_{-1} = \ln 2$	$\dot{v}_1 = \frac{1}{v_{-1}} \dot{v}_{-1} = \frac{1}{2}(1)$
$v_2 = v_{-1} \times v_0 = 2 \times 5$	$\dot{v}_2 = \dot{v}_{-1} \times v_0 + v_{-1} \times \dot{v}_0 = 1 \times 5 + 0 \times 2$
$v_3 = \sin(v_0)$	$\dot{v}_3 = \dot{v}_0 \times \cos(v_0)$
$v_4 = v_1 + v_2$	$\dot{v}_4 = \dot{v}_1 + \dot{v}_2$
$v_5 = v_4 - v_3$	$\dot{v}_5 = \dot{v}_4 - \dot{v}_3$
$y = v_5$	$\dot{y} = \dot{v}_5$

Forward primal trace and tangent trace for $\frac{\partial y}{\partial x_1}$ (II)

We now compute the derivative $\frac{\partial y}{\partial x_1}$ at $x_1 = 2, x_2 = 5$.

Forward primal trace	Forward tangent trace
$v_{-1} = x_1 = 2$	$\dot{v}_{-1} = \dot{x}_1 = 1$
$v_0 = x_2 = 5$	$\dot{v}_0 = \dot{x}_2 = 0$
$v_1 = \ln v_{-1} = \ln 2$	$\dot{v}_1 = \frac{1}{v_{-1}} \dot{v}_{-1} = \frac{1}{2}(1)$
$v_2 = v_{-1} \times v_0 = 2 \times 5$	$\dot{v}_2 = \dot{v}_{-1} \times v_0 + v_{-1} \times \dot{v}_0 = 1 \times 5 + 0 \times 2$
$v_3 = \sin(v_0) = \sin(5)$	$\dot{v}_3 = \dot{v}_0 \times \cos(v_0) = 0 \times \cos(5)$
$v_4 = v_1 + v_2$	$\dot{v}_4 = \dot{v}_1 + \dot{v}_2$
$v_5 = v_4 - v_3$	$\dot{v}_5 = \dot{v}_4 - \dot{v}_3$
$y = v_5$	$\dot{y} = \dot{v}_5$

Forward primal trace and tangent trace for $\frac{\partial y}{\partial x_1}$ (II)

We now compute the derivative $\frac{\partial y}{\partial x_1}$ at $x_1 = 2, x_2 = 5$.

Forward primal trace		Forward tangent trace	
$v_{-1} = x_1$	$= 2$	$\dot{v}_{-1} = \dot{x}_1$	$= 1$
$v_0 = x_2$	$= 5$	$\dot{v}_0 = \dot{x}_2$	$= 0$
$v_1 = \ln v_{-1}$	$= \ln 2$	$\dot{v}_1 = \frac{1}{v_{-1}} \dot{v}_{-1}$	$= \frac{1}{2}(1)$
$v_2 = v_{-1} \times v_0$	$= 2 \times 5$	$\dot{v}_2 = \dot{v}_{-1} \times v_0 + v_{-1} \times \dot{v}_0$	$= 1 \times 5 + 0 \times 2$
$v_3 = \sin(v_0)$	$= \sin(5)$	$\dot{v}_3 = \dot{v}_0 \times \cos(v_0)$	$= 0 \times \cos(5)$
$v_4 = v_1 + v_2$	$= 0.693 + 10$	$\dot{v}_4 = \dot{v}_1 + \dot{v}_2$	$= 0.5 + 5$
$v_5 = v_4 - v_3$		$\dot{v}_5 = \dot{v}_4 - \dot{v}_3$	
$y = v_5$		$\dot{y} = \dot{v}_5$	

Forward primal trace and tangent trace for $\frac{\partial y}{\partial x_1}$ (II)

We now compute the derivative $\frac{\partial y}{\partial x_1}$ at $x_1 = 2, x_2 = 5$.

Forward primal trace	Forward tangent trace
$v_{-1} = x_1 = 2$	$\dot{v}_{-1} = \dot{x}_1 = 1$
$v_0 = x_2 = 5$	$\dot{v}_0 = \dot{x}_2 = 0$
$v_1 = \ln v_{-1} = \ln 2$	$\dot{v}_1 = \frac{1}{v_{-1}} \dot{v}_{-1} = \frac{1}{2}(1)$
$v_2 = v_{-1} \times v_0 = 2 \times 5$	$\dot{v}_2 = \dot{v}_{-1} \times v_0 + v_0 \times \dot{v}_{-1} = 1 \times 5 + 0 \times 2$
$v_3 = \sin(v_0) = \sin(5)$	$\dot{v}_3 = \dot{v}_0 \times \cos(v_0) = 0 \times \cos(5)$
$v_4 = v_1 + v_2 = 0.693 + 10$	$\dot{v}_4 = \dot{v}_1 + \dot{v}_2 = 0.5 + 5$
$v_5 = v_4 - v_3 = 10.693 + 0.959$	$\dot{v}_5 = \dot{v}_4 - \dot{v}_3 = 5.5 - 0$
$y = v_5$	$\dot{y} = \dot{v}_5$

Forward primal trace and tangent trace for $\frac{\partial y}{\partial x_1}$ (II)

We now compute the derivative $\frac{\partial y}{\partial x_1}$ at $x_1 = 2, x_2 = 5$.

Forward primal trace		Forward tangent trace	
$v_{-1} = x_1$	$= 2$	$\dot{v}_{-1} = \dot{x}_1$	$= 1$
$v_0 = x_2$	$= 5$	$\dot{v}_0 = \dot{x}_2$	$= 0$
$v_1 = \ln v_{-1}$	$= \ln 2$	$\dot{v}_1 = \frac{1}{v_{-1}} \dot{v}_{-1}$	$= \frac{1}{2}(1)$
$v_2 = v_{-1} \times v_0$	$= 2 \times 5$	$\dot{v}_2 = \dot{v}_{-1} \times v_0 + v_0 \times \dot{v}_{-1}$	$= 1 \times 5 + 0 \times 2$
$v_3 = \sin(v_0)$	$= \sin(5)$	$\dot{v}_3 = \dot{v}_0 \times \cos(v_0)$	$= 0 \times \cos(5)$
$v_4 = v_1 + v_2$	$= 0.693 + 10$	$\dot{v}_4 = \dot{v}_1 + \dot{v}_2$	$= 0.5 + 5$
$v_5 = v_4 - v_3$	$= 10.693 + 0.959$	$\dot{v}_5 = \dot{v}_4 - \dot{v}_3$	$= 5.5 - 0$
$y = v_5$	$= 11.652$	$\dot{y} = \dot{v}_5$	$= 5.5$

Forward primal trace and tangent trace for $\frac{\partial y}{\partial x_1}$ (II)

We now compute the derivative $\frac{\partial y}{\partial x_1}$ at $x_1 = 2, x_2 = 5$.

which also means that to calculate
the derivatives wrt x_1

Forward primal trace	Forward tangent trace
$v_{-1} = x_1 = 2$	$\dot{v}_{-1} = \dot{x}_1 = 1$
$v_0 = x_2 = 5$	$\dot{v}_0 = \dot{x}_2 = 0$
$v_1 = \ln v_{-1} = \ln 2$	$\dot{v}_1 = \frac{1}{v_{-1}} \dot{v}_{-1} = \frac{1}{2}(1)$
$v_2 = v_{-1} \times v_0 = 2 \times 5$	$\dot{v}_2 = \dot{v}_{-1} \times v_0 + v_0 \times \dot{v}_{-1} = 1 \times 5 + 0 \times 2$
$v_3 = \sin(v_0) = \sin(5)$	$\dot{v}_3 = \dot{v}_0 \times \cos(v_0) = 0 \times \cos(5)$
$v_4 = v_1 + v_2 = 0.693 + 10$	$\dot{v}_4 = \dot{v}_1 + \dot{v}_2 = 0.5 + 5$
$v_5 = v_4 - v_3 = 10.693 + 0.959$	$\dot{v}_5 = \dot{v}_4 - \dot{v}_3 = 5.5 - 0$
$y = v_5 = 11.652$	$\dot{y} = \dot{v}_5 = 5.5$

If we want to compute $\frac{\partial y}{\partial x_2}$ instead, we set $\dot{v}_{-1} = 0$ and $\dot{v}_0 = 1$.

Generalisation to the Jacobian of a function

- Let $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ be a function with n independent variables x_i and m dependent variables y_j .
- The derivatives in the Jacobian, $\frac{\partial y_j}{\partial x_i}$, are computed by making $\dot{x}_i = 1$ initially in the forward pass and all the other derivatives $\dot{x}_k = 0$ for $k \neq i$.
- The values of the derivatives at $\mathbf{x} = \mathbf{a}$,

$$\dot{y}_j = \left. \frac{\partial y_j}{\partial x_i} \right|_{\mathbf{x}=\mathbf{a}}.$$

are obtained by a forward pass of AD.

- Notice that for a specific x_i , we can compute all the derivatives $\frac{\partial y_j}{\partial x_i}$ for $j = 1, \dots, m$, which corresponds to the column i in the Jacobian.
- To compute the whole Jacobian, we need n forward passes, one per input variable.

Complexity

Expensive when computing
choose reverse mode of AD $n \gg m$.

- AD with forward mode is efficient for functions like $\mathbf{f} : \mathbb{R} \rightarrow \mathbb{R}^m$.
- The reason, as we saw before, is because we can compute all the derivatives $\frac{\partial y_j}{\partial x}$ for $j = 1, \dots, m$ in one pass.
- In the other extreme, $f : \mathbb{R}^n \rightarrow \mathbb{R}$, it needs n forward passes and it can become **computationally expensive** when n is large.
- In general, when $n \gg m$, the reverse mode of AD is preferred.

Contents

Derivatives and ways to compute them

AD modes

- Forward mode
- Reverse mode

Implementations

Backpropagate

- AD in reverse mode propagates derivatives backwards from a given output.
- It is done by computing intermediate variables for v_i known as *adjoints*,

$$\bar{v}_i = \frac{\partial y_j}{\partial v_i}, \quad v_i \text{ bar}$$

representing the sensitivity of output y_j to input v_i .

- AD in reverse mode uses two-phases
 - a *forward* step to compute the variables v_i and to book-keep dependencies in the computational graph.
 - a *backward* or *reverse* step, in which the adjoints are used to compute the derivatives, starting from the outputs and going back to the inputs.

Example (I)

- Let us go back to the example we saw before

$$y = f(x_1, x_2) = \ln(x_1) + x_1 x_2 - \sin(x_2),$$

and focus on v_0 ($v_0 = x_2$).

- We want to compute the adjoint $\bar{v}_0 = \frac{\partial y}{\partial v_0}$, this is, how the change in v_0 affects the output y .
- From the computational graph, we see that $\underline{v_0}$ affects y through v_2 and v_3 ,

$$\begin{aligned}\frac{\partial y}{\partial v_0} &= \frac{\partial y}{\partial z} + \frac{\partial y}{\partial v_4} \\&= \frac{\partial y}{\partial v_5} = \frac{\partial y}{\partial v_5} \cdot \frac{\partial v_5}{\partial v_0} \\&\frac{\partial y}{\partial p} = \frac{\partial y}{\partial p} \cdot \frac{\partial p}{\partial z} \cdot \frac{\partial z}{\partial v_0}\end{aligned}$$
$$\begin{aligned}&v_5 \rightarrow y \\&\frac{\partial y}{\partial v_5} = \frac{\partial y}{\partial v_5} \cdot \frac{\partial v_5}{\partial v_0}\end{aligned}$$

Example (II) $\frac{dy}{dv_0} =$

$$\frac{dy}{dv_2} \cdot \frac{dv_2}{dv_0}$$

- So the contribution of v_0 to y is given as

$$\frac{\partial y}{\partial v_0} = \frac{\partial y}{\partial v_2} \frac{\partial v_2}{\partial v_0} + \frac{\partial y}{\partial v_3} \frac{\partial v_3}{\partial v_0}. \quad \text{direct way}$$

Total Gradient

By definition $\frac{\partial y}{\partial v_2} = \bar{v}_2$ and $\frac{\partial y}{\partial v_3} = \bar{v}_3$, so we can write the expression above as

adjoint for v_2 :

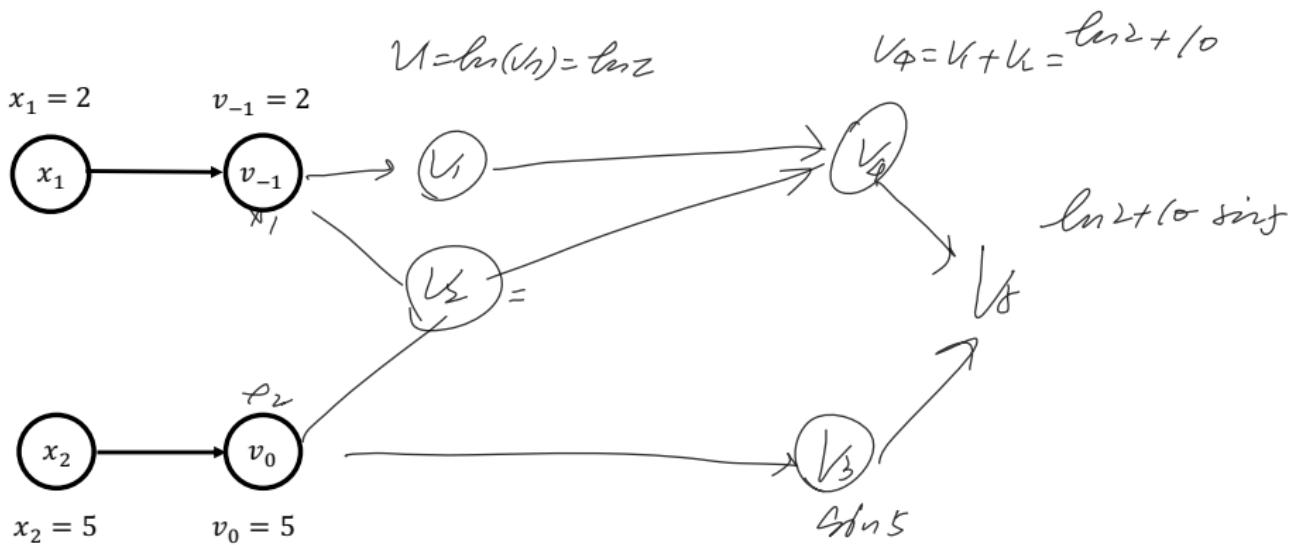
$$\frac{dy}{dx} = \bar{v}_2 \cdot \frac{dy}{dv_2} = \bar{v}_2 \cdot \frac{dy}{dv_2} + \bar{v}_3 \cdot \frac{dy}{dv_3} = \bar{v}_2 + \bar{v}_3 \frac{\partial v_3}{\partial v_0}.$$

$$\bar{v}_0 = \bar{v}_2 \cdot \frac{dv_2}{dv_0} + \bar{v}_3 \cdot \frac{dv_3}{dv_0}$$

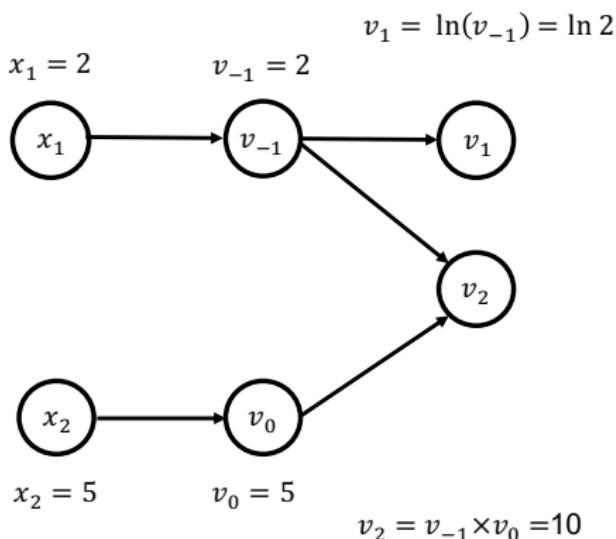
- After the forward pass to compute v_i , the reverse pass computes the adjoints, starting with $\bar{v}_5 = \bar{y} = \frac{\partial y}{\partial y} = 1$, and computing $\frac{\partial y}{\partial x_1} = \bar{x}_1$ and $\frac{\partial y}{\partial x_2} = \bar{x}_2$ at the end.

reverse mode

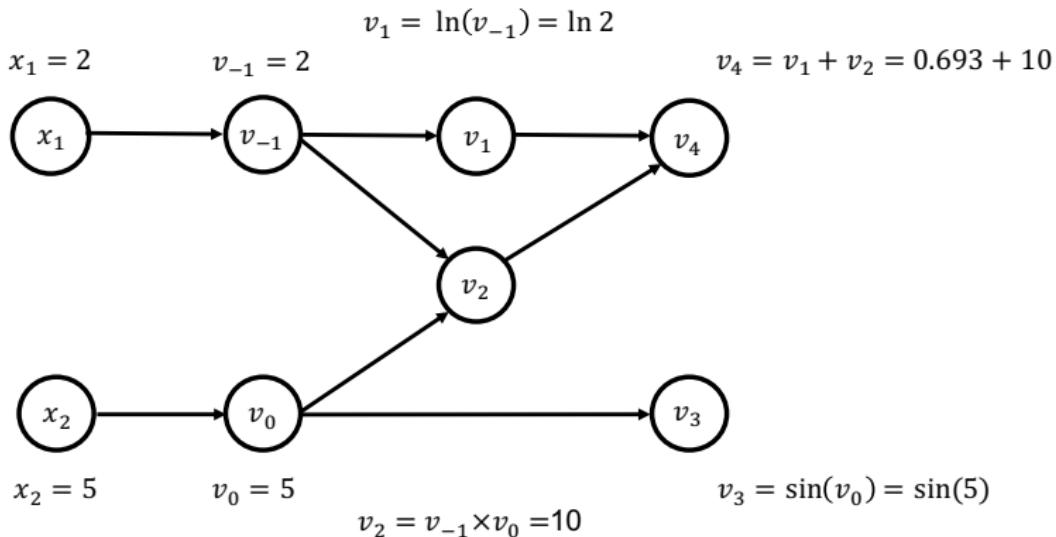
Forward primal trace (forward pass)



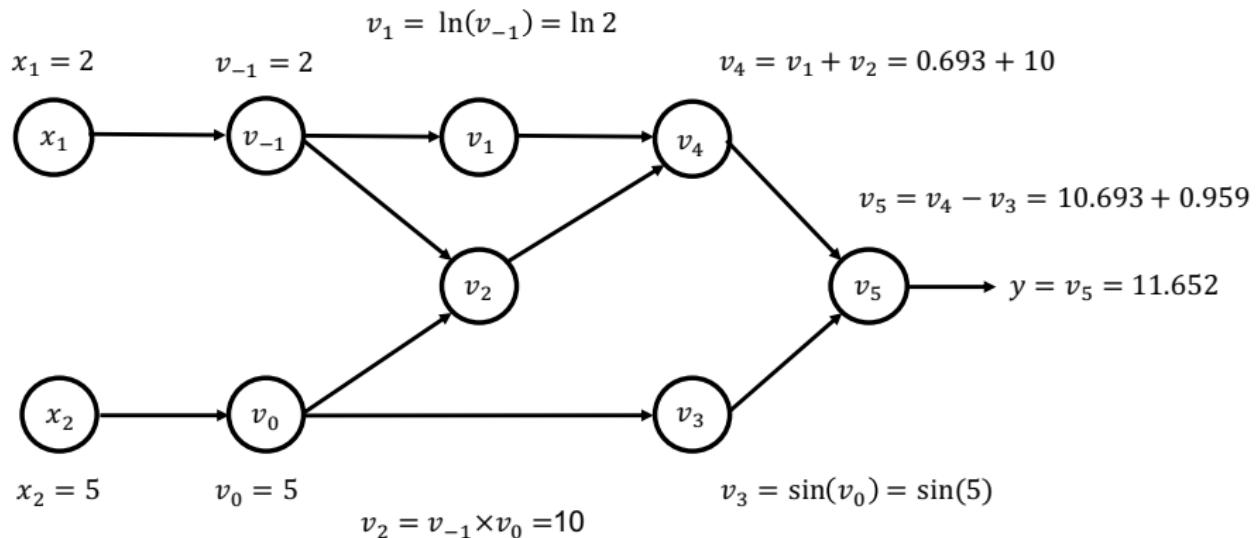
Forward primal trace (forward pass)



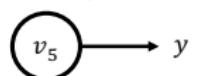
Forward primal trace (forward pass)



Forward primal trace (forward pass)



Reverse adjoint (derivative) trace (reverse pass)

$$\bar{v}_5 = \bar{y} = 1$$


A circular node with the label v_5 inside. An arrow originates from the right side of the circle and points to the right, ending at a small tick mark. To the right of the arrow tip, the label y is written.

Reverse adjoint (derivative) trace (reverse pass)

$$\bar{v}_4 = \frac{dy}{dv_4} = \frac{dy}{dv_5} \cdot \frac{dv_5}{dv_4} = \bar{v}_5 \cdot \frac{dv_5}{dv_4}$$

v_4

$$\bar{v}_5 = \bar{y} = 1$$

v_5

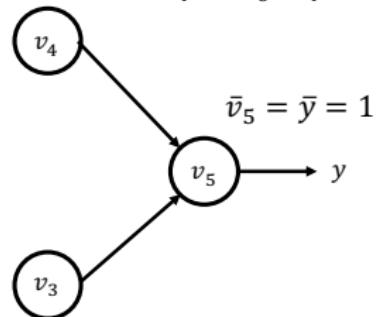
y

v_3

$$\bar{v}_3 = \frac{dy}{dv_3} = \frac{dy}{dv_5} \cdot \frac{dv_5}{dv_3} = \bar{v}_5 \cdot \frac{dv_5}{dv_3}$$

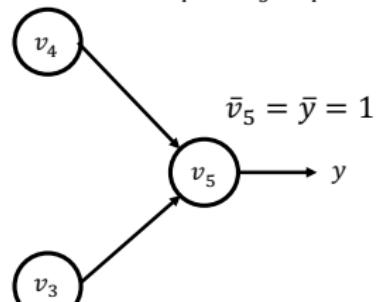
Reverse adjoint (derivative) trace (reverse pass)

$$\bar{v}_4 = \frac{\partial y}{\partial v_4} = \frac{\partial y}{\partial v_5} \frac{\partial v_5}{\partial v_4} = \bar{v}_5 \frac{\partial v_5}{\partial v_4}$$



Reverse adjoint (derivative) trace (reverse pass)

$$\bar{v}_4 = \frac{\partial y}{\partial v_4} = \frac{\partial y}{\partial v_5} \frac{\partial v_5}{\partial v_4} = \bar{v}_5 \frac{\partial v_5}{\partial v_4}$$

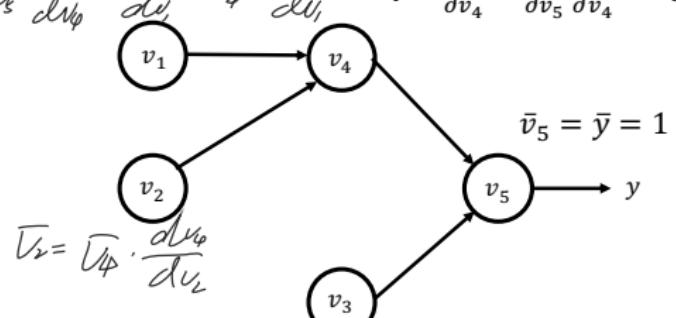


$$\bar{v}_3 = \frac{\partial y}{\partial v_3} = \frac{\partial y}{\partial v_5} \frac{\partial v_5}{\partial v_3} = \bar{v}_5 \frac{\partial v_5}{\partial v_3}$$

Reverse adjoint (derivative) trace (reverse pass)

$$\bar{v}_1 = \frac{\partial y}{\partial v_1} = \frac{\partial y}{\partial v_5} \frac{\partial v_5}{\partial v_4} \frac{\partial v_4}{\partial v_1} = \bar{v}_4 \frac{\partial v_4}{\partial v_1}$$

$$\bar{v}_1 = \frac{dy}{du} = \frac{dy}{ds} \cdot \frac{ds}{dp} \cdot \frac{de}{dv_i} = \bar{v}_4 \cdot \frac{ds}{dv_i} \cdot \frac{dv_i}{du} = \bar{v}_4 \cdot \frac{dv_i}{du}, \quad \bar{v}_4 = \frac{\partial y}{\partial v_4} = \frac{\partial y}{\partial v_5} \frac{\partial v_5}{\partial v_4} = \bar{v}_5 \frac{\partial v_5}{\partial v_4}$$



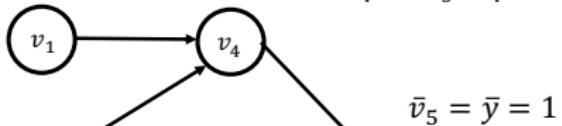
$$\bar{v}_2 = \bar{v}_4 \cdot \frac{dv_4}{dv_i}$$

$$\bar{v}_3 = \frac{\partial y}{\partial v_3} = \frac{\partial y}{\partial v_5} \frac{\partial v_5}{\partial v_3} = \bar{v}_5 \frac{\partial v_5}{\partial v_3}$$

Reverse adjoint (derivative) trace (reverse pass)

$$\bar{v}_1 = \frac{\partial y}{\partial v_1} = \frac{\partial y}{\partial v_5} \frac{\partial v_5}{\partial v_4} \frac{\partial v_4}{\partial v_1} = \bar{v}_4 \frac{\partial v_4}{\partial v_1}$$

$$\bar{v}_4 = \frac{\partial y}{\partial v_4} = \frac{\partial y}{\partial v_5} \frac{\partial v_5}{\partial v_4} = \bar{v}_5 \frac{\partial v_5}{\partial v_4}$$



$$\bar{v}_3 = \frac{\partial y}{\partial v_3} = \frac{\partial y}{\partial v_5} \frac{\partial v_5}{\partial v_3} = \bar{v}_5 \frac{\partial v_5}{\partial v_3}$$

$$\bar{v}_2 = \frac{\partial y}{\partial v_2} = \frac{\partial y}{\partial v_5} \frac{\partial v_5}{\partial v_4} \frac{\partial v_4}{\partial v_2} = \bar{v}_4 \frac{\partial v_4}{\partial v_2}$$

Reverse adjoint (derivative) trace (reverse pass)

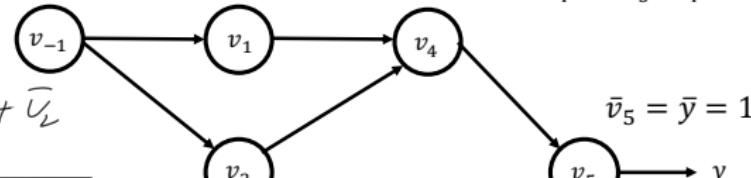
$$\bar{v}_i = \frac{\partial y}{\partial v_i}$$

$$\bar{v}_1 = \frac{\partial y}{\partial v_1} = \frac{\partial y}{\partial v_5} \frac{\partial v_5}{\partial v_4} \frac{\partial v_4}{\partial v_1} = \bar{v}_4 \frac{\partial v_4}{\partial v_1}$$

$$\bar{v}_{-1} = \frac{\partial y}{\partial v_{-1}} = \frac{\partial y}{\partial v_1} \frac{\partial v_1}{\partial v_{-1}} + \frac{\partial y}{\partial v_2} \frac{\partial v_2}{\partial v_{-1}} = \bar{v}_1 \frac{\partial v_1}{\partial v_{-1}} + \bar{v}_2 \frac{\partial v_2}{\partial v_{-1}}$$

$$\bar{v}_4 = \frac{\partial y}{\partial v_4} = \frac{\partial y}{\partial v_5} \frac{\partial v_5}{\partial v_4} = \bar{v}_5 \frac{\partial v_5}{\partial v_4}$$

$$\begin{aligned} \bar{v}_i &= \frac{\partial y}{\partial y} \cdot \frac{\partial y}{\partial v_i} \left(\frac{\partial v_1}{\partial v_i} + \frac{\partial v_2}{\partial v_i} \right) \\ &= \bar{v}_4 \cdot \frac{\partial v_1}{\partial v_i} + \bar{v}_5 \cdot \frac{\partial v_2}{\partial v_i} = \bar{v}_i + \bar{v}_i \end{aligned}$$



$$\bar{v}_0 = \frac{\partial y}{\partial v_0} \cdot \frac{\partial v_1}{\partial v_0} \cdot \frac{\partial v_2}{\partial v_0} + \bar{v}_2 \cdot \frac{\partial v_1}{\partial v_0}$$

different places

$$\bar{v}_3 = \frac{\partial y}{\partial v_3} = \frac{\partial y}{\partial v_5} \frac{\partial v_5}{\partial v_3} = \bar{v}_5 \frac{\partial v_5}{\partial v_3}$$

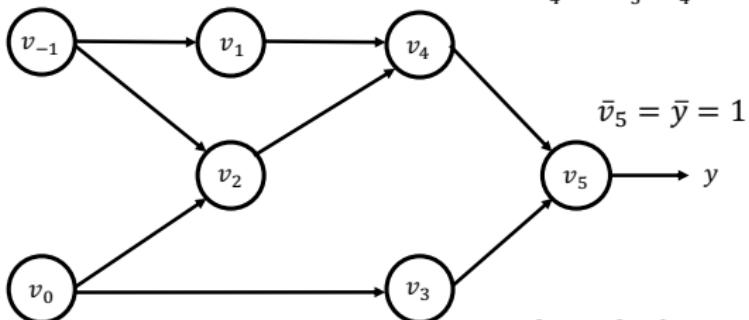
$$\bar{v}_2 = \frac{\partial y}{\partial v_2} = \frac{\partial y}{\partial v_5} \frac{\partial v_5}{\partial v_4} \frac{\partial v_4}{\partial v_2} = \bar{v}_4 \frac{\partial v_4}{\partial v_2}$$

$$\begin{aligned} \frac{\partial v_1}{\partial v_3} \cdot \frac{\partial v_2}{\partial v_3} \cdot \frac{\partial v_1}{\partial v_0} &= \bar{v}_3 \cdot \frac{\partial v_2}{\partial v_0} \\ &= \bar{v}_2 \cdot \frac{\partial v_1}{\partial v_0} + \bar{v}_3 \cdot \frac{\partial v_2}{\partial v_0} \end{aligned}$$

Reverse adjoint (derivative) trace (reverse pass)

$$\bar{v}_1 = \frac{\partial y}{\partial v_1} = \frac{\partial y}{\partial v_5} \frac{\partial v_5}{\partial v_4} \frac{\partial v_4}{\partial v_1} = \bar{v}_4 \frac{\partial v_4}{\partial v_1}$$

$$\bar{v}_{-1} = \frac{\partial y}{\partial v_{-1}} = \frac{\partial y}{\partial v_1} \frac{\partial v_1}{\partial v_{-1}} + \frac{\partial y}{\partial v_2} \frac{\partial v_2}{\partial v_{-1}} = \bar{v}_1 \frac{\partial v_1}{\partial v_{-1}} + \bar{v}_2 \frac{\partial v_2}{\partial v_{-1}} \quad \bar{v}_4 = \frac{\partial y}{\partial v_4} = \frac{\partial y}{\partial v_5} \frac{\partial v_5}{\partial v_4} = \bar{v}_5 \frac{\partial v_5}{\partial v_4}$$



$$\bar{v}_0 = \frac{\partial y}{\partial v_0} = \frac{\partial y}{\partial v_2} \frac{\partial v_2}{\partial v_0} + \frac{\partial y}{\partial v_3} \frac{\partial v_3}{\partial v_0} = \bar{v}_2 \frac{\partial v_2}{\partial v_0} + \bar{v}_3 \frac{\partial v_3}{\partial v_0}$$

$$\bar{v}_3 = \frac{\partial y}{\partial v_3} = \frac{\partial y}{\partial v_5} \frac{\partial v_5}{\partial v_3} = \bar{v}_5 \frac{\partial v_5}{\partial v_3}$$

$$\bar{v}_2 = \frac{\partial y}{\partial v_2} = \frac{\partial y}{\partial v_5} \frac{\partial v_5}{\partial v_4} \frac{\partial v_4}{\partial v_2} = \bar{v}_4 \frac{\partial v_4}{\partial v_2}$$

Reverse adjoint (derivative) trace (reverse pass)

$$\bar{v}_1 = \frac{\partial y}{\partial v_1} = \frac{\partial y}{\partial v_5} \frac{\partial v_5}{\partial v_4} \frac{\partial v_4}{\partial v_1} = \bar{v}_4 \frac{\partial v_4}{\partial v_1}$$

$$\bar{v}_1 = \frac{dy}{dx} = \frac{\bar{v}_1}{\bar{v}_2} = \frac{dy}{dv_1} = \frac{\partial y}{\partial v_1} \frac{\partial v_1}{\partial v_{-1}} + \frac{\partial y}{\partial v_2} \frac{\partial v_2}{\partial v_{-1}} = \bar{v}_1 \frac{\partial v_1}{\partial v_{-1}} + \bar{v}_2 \frac{\partial v_2}{\partial v_{-1}}$$

$$\bar{x}_1 = \frac{\partial y}{\partial x_1} = \frac{\partial y}{\partial v_{-1}} \frac{\partial v_{-1}}{\partial x_1} = \bar{v}_{-1} \frac{\partial v_{-1}}{\partial x_1}$$

$$\bar{x}_2 = \frac{\partial y}{\partial x_2} = \frac{\partial y}{\partial v_0} \frac{\partial v_0}{\partial x_2} = \cancel{\bar{v}_0} \frac{\partial v_0}{\partial x_2}$$

$$\bar{v}_0 = \frac{\partial y}{\partial v_0} = \frac{\partial y}{\partial v_2} \frac{\partial v_2}{\partial v_0} + \frac{\partial y}{\partial v_3} \frac{\partial v_3}{\partial v_0} = \bar{v}_2 \frac{\partial v_2}{\partial v_0} + \bar{v}_3 \frac{\partial v_3}{\partial v_0}$$

$$\bar{v}_3 = \frac{\partial y}{\partial v_3} = \frac{\partial y}{\partial v_5} \frac{\partial v_5}{\partial v_3} = \bar{v}_5 \frac{\partial v_5}{\partial v_3}$$

$$\bar{v}_2 = \frac{\partial y}{\partial v_2} = \frac{\partial y}{\partial v_5} \frac{\partial v_5}{\partial v_4} \frac{\partial v_4}{\partial v_2} = \bar{v}_4 \frac{\partial v_4}{\partial v_2}$$

$$\bar{v}_i = \bar{l}_i \cdot \frac{dv_i}{dx_i}$$

Numerical evaluation of the reverse adjoint trace

reverse adjoint trace

why? $\frac{d\bar{v}_5}{dV_4} = 1?$

$$\bar{v}_4 = \bar{v}_5 \frac{\partial v_5}{\partial v_4} = \bar{v}_5 \times 1 = 1$$

$$v_4 = V_1 + V_2$$

$$\bar{v}_5 = \bar{y} = 1 \quad \frac{dy}{dV_4}$$

$$v_5 = V_1 + V_2 - V_3$$

$$v_3 = \sin(xy)$$

$$\bar{v}_3 = \bar{v}_5 \frac{\partial v_5}{\partial v_3} = \bar{v}_5 \times (-1) = -1$$

why $\frac{d\bar{v}_5}{dV_3} = -1$

If you don't know something

just follow more steps

Numerical evaluation of the reverse adjoint trace

$$v_1 = \ln(x_1)$$

$$v_2 = x_1 \cdot x_2$$

$$v_3 = \sin(x_2)$$

$$v_4 = v_1 + v_2$$

$$v_5 = v_4 - v_3$$

$$v_6 = v_1 + v_3$$

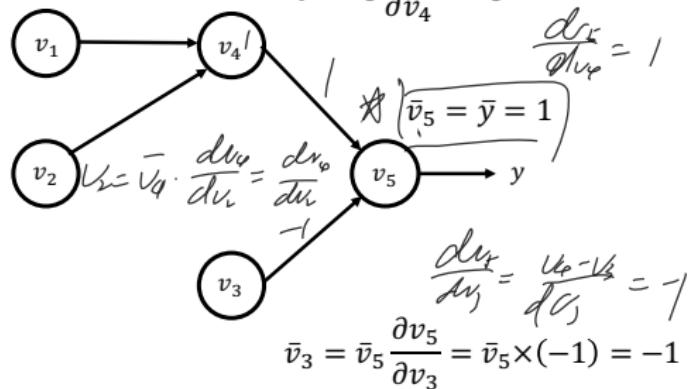
$$\frac{dv_6}{dv_1} = \left(\frac{v_1 + v_3}{v_1} \right) \cancel{\frac{dx_1}{dv_1}} = 1$$

$$\bar{v}_1 = \bar{v}_4 \frac{\partial v_4}{\partial v_1} = \bar{v}_4 \times 1 = 1$$

$$l_5 =$$

$$\bar{v}_4 = \bar{v}_5 \frac{\partial v_5}{\partial v_4} = \bar{v}_5 \times 1 = 1$$

$$\cancel{\frac{dx_5}{dv_4}} = 1$$



$$\bar{v}_2 = \bar{v}_4 \frac{\partial v_4}{\partial v_2} = \bar{v}_4 \times 1 = 1$$

$$\frac{dy}{dv_3} = \frac{v_5 - v_4}{v_5} = -1$$

$$\bar{v}_3 = \bar{v}_5 \frac{\partial v_5}{\partial v_3} = \bar{v}_5 \times (-1) = -1$$

Numerical evaluation of the reverse adjoint trace

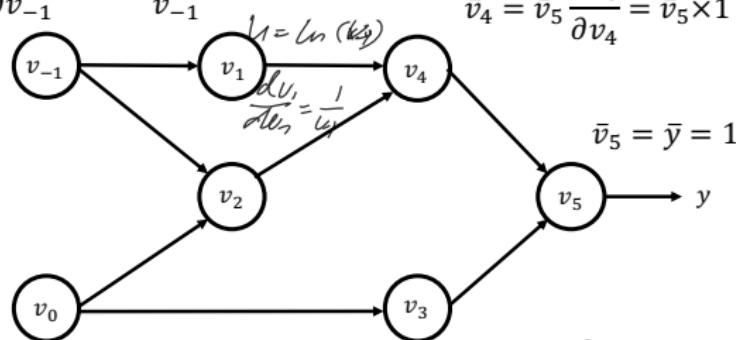
$$U_r = U_1 \cdot U_0$$

$$\frac{dU_r}{dV_1} = U_1$$

$$\bar{v}_1 = \bar{v}_4 \frac{\partial v_4}{\partial v_1} = \bar{v}_4 \times 1 = 1$$

$$\bar{v}_{-1} = \bar{v}_1 \frac{\partial v_1}{\partial v_{-1}} + \bar{v}_2 \frac{\partial v_2}{\partial v_{-1}} = \bar{v}_1 \times \frac{1}{v_{-1}} + \bar{v}_2 \times v_0 = 5.5$$

$$\bar{v}_4 = \bar{v}_5 \frac{\partial v_5}{\partial v_4} = \bar{v}_5 \times 1 = 1$$



$$\bar{v}_0 = \bar{v}_2 \frac{\partial v_2}{\partial v_0} + \bar{v}_3 \frac{\partial v_3}{\partial v_0} = \bar{v}_2 \times v_{-1} + \bar{v}_3 \cos(v_0) = 1.716$$

$$\bar{v}_3 = \bar{v}_5 \frac{\partial v_5}{\partial v_3} = \bar{v}_5 \times (-1) = -1$$

$$\bar{v}_2 = \bar{v}_4 \frac{\partial v_4}{\partial v_2} = \bar{v}_4 \times 1 = 1$$

Numerical evaluation of the reverse adjoint trace

$$\bar{v}_1 = \bar{v}_4 \frac{\partial v_4}{\partial v_1} = \bar{v}_4 \times 1 = 1$$

$$\bar{v}_{-1} = \bar{v}_1 \frac{\partial v_1}{\partial v_{-1}} + \bar{v}_2 \frac{\partial v_2}{\partial v_{-1}} = \bar{v}_1 \times \frac{1}{v_{-1}} + \bar{v}_2 \times v_0 = 5.5$$

$$\bar{v}_4 = \bar{v}_5 \frac{\partial v_5}{\partial v_4} = \bar{v}_5 \times 1 = 1$$

x_1

v_{-1}

v_1

v_4

v_2

v_3

v_5

y

$$\bar{x}_1 = \bar{v}_{-1} \frac{\partial v_{-1}}{\partial x_1} = 5.5 \times 1 = 5.5$$

$$\bar{v}_5 = \bar{y} = 1$$

$$\bar{x}_2 = \bar{v}_0 \frac{\partial v_0}{\partial x_2} = 1.716 \times 1 = 1.716$$

$$y$$

x_2

v_0

v_1

v_2

v_3

v_4

v_5

$$\bar{v}_0 = \bar{v}_2 \frac{\partial v_2}{\partial v_0} + \bar{v}_3 \frac{\partial v_3}{\partial v_0} = \bar{v}_2 \times v_{-1} + \bar{v}_3 \cos(v_0) = 1.716$$

$$\bar{v}_3 = \bar{v}_5 \frac{\partial v_5}{\partial v_3} = \bar{v}_5 \times (-1) = -1$$

$$\bar{v}_2 = \bar{v}_4 \frac{\partial v_4}{\partial v_2} = \bar{v}_4 \times 1 = 1$$

Complexity

- Reverse mode AD performs better when $n \gg m$.
why? why it performs better.
- The downside is the cost of increased storage, since we need to save intermediate values for v_i in the evaluation trace.

Because it starts at the end

For forward trace, do not need to save the intermediate values. → just calculate is good enough to get \mathbf{u} and \mathbf{v} .

Reverse mode AD and backpropagation

- Reverse mode AD is the algorithm used to train neural networks and deep learning models.
- To train a neural network model, we optimise an objective function, $E(\mathbf{w}) : \mathbb{R}^n \rightarrow \mathbb{R}^m$ that usually depends on a high-dimensional input vector of parameters $\mathbf{w} \in \mathbb{R}^n$, with $n \gg m$.
*really high dimensional input vector → a relatively small number of results.
we consider several factors (A-lot-of)*
- In the machine learning community, reverse mode AD goes by the name of **backpropagation**, which you will see again in the session on neural networks.

Contents

Derivatives and ways to compute them

AD modes

- Forward mode
- Reverse mode

Implementations

AD implementations

Table 5: Survey of AD implementations. Tools developed primarily for machine learning are highlighted in bold.

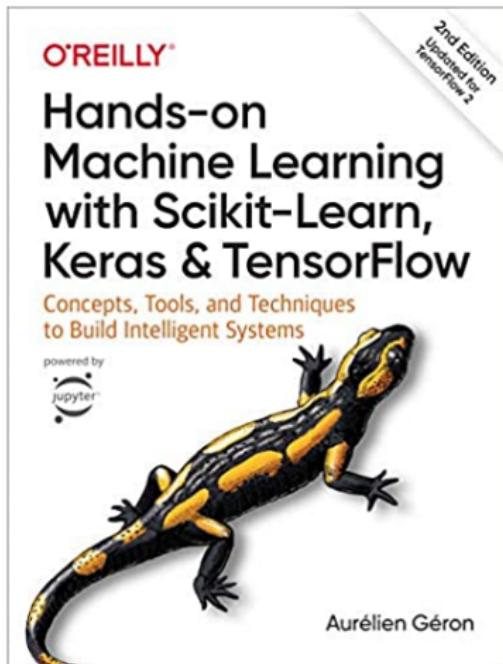
Language	Tool	Type	Mode	Institution / Project	Reference	URL
AMPL	AMPL	INT	F, R	Bell Laboratories	Fourer et al. (2002)	http://www.ampl.com/
C, C++	ADIC	ST	F, R	Argonne National Laboratory	Bischof et al. (1997)	http://www.mcs.anl.gov/research/projects/adic/
	ADOL-C	OO	F, R	Computational Infrastructure for Operations Research	Walther and Griewank (2012)	https://projects.coin-or.org/ADOL-C
C++	Ceres Solver	LIB	F	Google		http://ceres-solver.org/
	CppAD	OO	F, R	Computational Infrastructure for Operations Research	Bell and Burke (2008)	http://www.coin-or.org/CppAD/
	FABDAD++	OO	F, R	Technical University of Denmark	Bendtsen and Stauning (1996)	http://www.fabdad.com/fabdad.html
	MxPyptk	OO	F	Fermi National Accelerator Laboratory	Ostiguy and Michelotti (2007)	
C#	AutoDiff	LIB	R	George Mason Univ., Dept. of Computer Science	Shtof et al. (2013)	http://autodiff.codeplex.com/
F#, C#	DiffSharp	OO	F, R	Maynooth University, Microsoft Research Cambridge	Baydin et al. (2016a)	http://diffsharp.github.io
Fortran	ADIFOR	ST	F, R	Argonne National Laboratory	Bischof et al. (1996)	http://www.mcs.anl.gov/research/projects/adifor/
	NAGWare	COM	F, R	Numerical Algorithms Group	Naumann and Riehme (2005)	http://www.nag.co.uk/nagware/Research/ad_overview.asp
	TAMC	ST	R	Max Planck Institute for Meteorology	Giering and Kaminski (1998)	http://autodiff.com/tamc/
Fortran, C	COSY	INT	F	Michigan State Univ., Biomedical and Physical Sci.	Berz et al. (1996)	http://www.bt.pa.msu.edu/index_cosy.htm
	Tapenade	ST	F, R	INRIA Sophia-Antipolis	Hascoët and Pascual (2013)	http://www-sop.inria.fr/tropics/tapenade.html
Haskell	ad	OO	F, R	Haskell package		http://hackage.haskell.org/package/ad
Java	ADJac	ST	F, R	University Politehnica of Bucharest	Slusanschi and Dumitrel (2016)	http://adjac.cs.pub.ro
	Deriva	LIB	R	Java & Clojure library		https://github.com/lambda/Deriva
Julia	JuliaDiff	OO	F, R	Julia packages	Revels et al. (2016a)	http://www.juliadiff.org/
Lua	torch-autograd	OO	R	Twitter Cortex		https://github.com/twitter/torch-autograd
MATLAB	ADiMat	ST	F, R	Technical University of Darmstadt, Scientific Comp.	Willkomm and Vehreschild (2013)	http://adimat.sc.informatik.tu-darmstadt.de/
	INTLab	OO	F	Hamburg Univ. of Technology, Inst. for Reliable Comp.	Rump (1999)	http://www.ti3.tu-harburg.de/rump/intlab/
	TOMLAB/MAD	OO	F	Cranfield University & Tomlab Optimization Inc.	Forth (2006)	http://tomlab.biz/products/mad
Python	ad	OO	R	Python package		https://pypi.python.org/pypi/ad
	autograd	OO	F, R	Harvard Intelligent Probabilistic Systems Group	Maclaurin (2016)	https://github.com/HIPS/autograd
	Chainer	OO	R	Preferred Networks	Tokui et al. (2015)	https://chainer.org/
	PyTorch	OO	R	PyTorch core team	Paszke et al. (2017)	http://pytorch.org/
	Tangent	ST	F, R	Google Brain	van Merriënboer et al. (2017)	https://github.com/google/tangent
Scheme	R6RS-AD	OO	F, R	Purdue Univ., School of Electrical and Computer Eng.		https://github.com/qobi/R6RS-AD
	Scmutils	OO	F	MIT Computer Science and Artificial Intelligence Lab.	Sussman and Wisdom (2001)	http://groups.csail.mit.edu/mac/users/gjs/6946/refman.txt
	Stalingrad	COM	F, R	Purdue Univ., School of Electrical and Computer Eng.	Pearlmutter and Siskind (2008)	http://www.bcl.hamilton.ie/~qobi/stalingrad/

F: Forward, R: Reverse; COM: Compiler, INT: Interpreter, LIB: Library, OO: Operator overloading, ST: Source transformation

Two popular implementations in the ML community



References



Appendix D of “Hands-On Machine Learning with Scikit-Learn, Keras and TensorFlow”

References

Journal of Machine Learning Research 18 (2018) 1-43

Submitted 8/17; Published 4/18

Automatic Differentiation in Machine Learning: a Survey

Atilim Güneş Baydin

*Department of Engineering Science
University of Oxford
Oxford OX1 3PJ, United Kingdom*

GUNES@ROBOTS.OX.AC.UK

Barak A. Pearlmutter

*Department of Computer Science
National University of Ireland Maynooth
Maynooth, Co. Kildare, Ireland*

BARAK@PEARLMUTTER.NET

Alexey Andreyevich Radul

*Department of Brain and Cognitive Sciences
Massachusetts Institute of Technology
Cambridge, MA 02139, United States*

AXCH@MIT.EDU

Jeffrey Mark Siskind

*School of Electrical and Computer Engineering
Purdue University
West Lafayette, IN 47907, United States*

QOBI@PURDUE.EDU