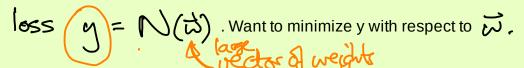
Machine learning, reverse-mode AD aka back-propagation

# Reverse mode AD cont'd

The above process is called "forward-mode" AD. It is accurate to machine epsilon, and easy to program. But it is not optimally efficient.

Machine learning - training a neural network - involves minimizing a real-valued function called the "loss", which is a measure of how poorly the network is performing the desired task.

The network is characterized by a large set of parameters called "weights",



The standard method is to move repeatedly in the direction of steepest descent. gradient descent

That direction is -

Because of the huge dimension of, we need to be as efficient as possible (as well as accurate).

"Reverse-mode AD" or "back-propagation" is more efficient than forward-mode AD.

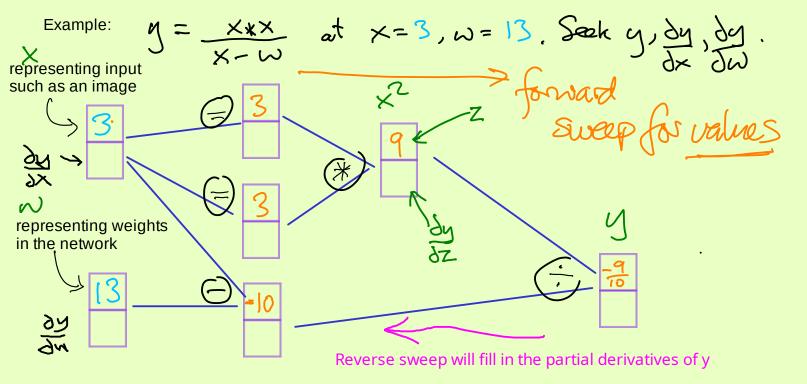
In the diagram below,

the upper slots are for the \*values\* of initial, intermediate, and final quantities in the calculation, — the lower slots are for the \*derivative of the output value y with respect to the quantity\*.

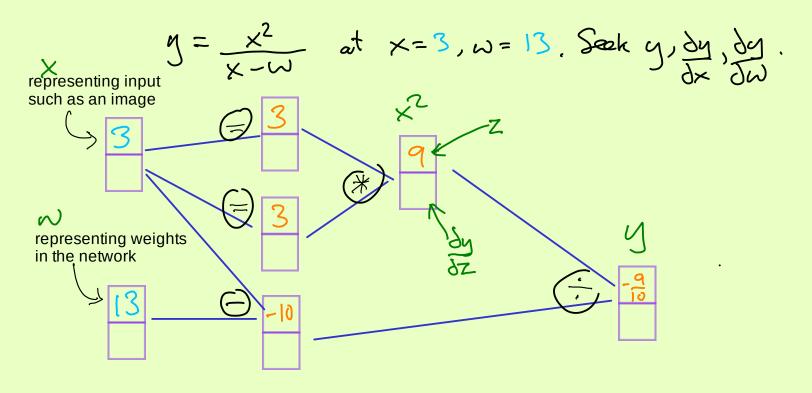
Our goal is to fill in all the slots in two sweeps:

- (i) a forward sweep to fill in all the \*values\*
- (ii) a reverse sweep (back propagation) to fill in all the partial derivatives,

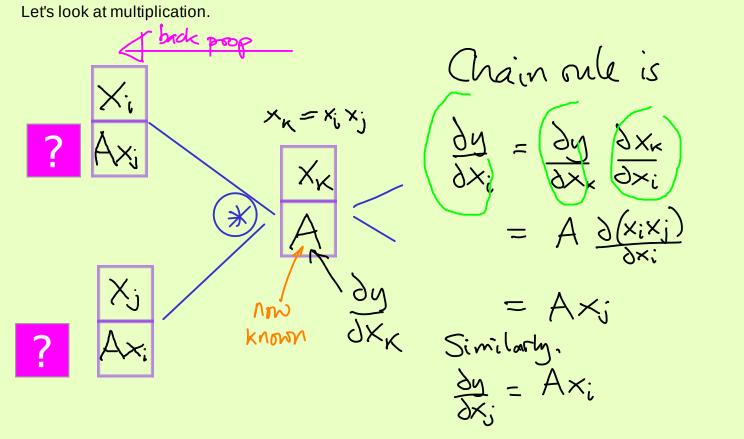
The last numbers we compute in the reverse sweep will be the desired components of the gradient of y.



Forward sweep to fill in the values:

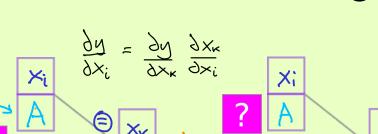


For the reverse sweep to fill in the partial derivatives of y, we need to develop some rules ...



### Exercise for you:

Find the back prop rules for the operations of (=), (-

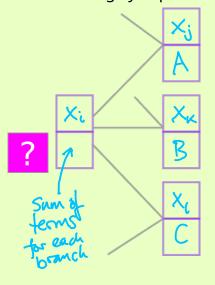


$$\frac{\partial u}{\partial x_i} = \frac{\partial u}{\partial x_k} \frac{\partial x_i}{\partial x_i} = \frac{\partial u}{\partial x_k} \frac{\partial x_i}{\partial x_i}$$

? 
$$\frac{\partial x_{i}}{\partial x_{j}} = \frac{\partial x_{k}}{\partial x_{k}} \frac{\partial x_{k}}{\partial x_{i}} = \frac{A \cdot (-i)}{A \cdot (-i)}$$

$$\frac{\partial x_{i}}{\partial x_{i}} = \frac{\partial x_{i}}{\partial x_{i}} \frac{\partial x_{i}}{\partial x_{i}} = \frac{\partial x_$$

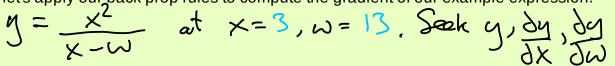
## Branching - y depends on **X**ivia multiple paths

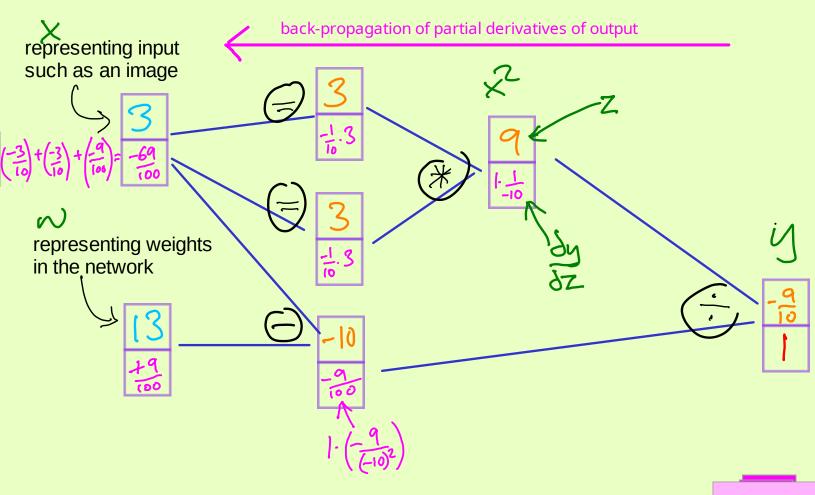


$$\frac{\partial g}{\partial x_i} = \frac{\partial g}{\partial x_j} \frac{\partial x_i}{\partial x_i} + \frac{\partial g}{\partial x_k} \frac{\partial x_k}{\partial x_i} + \frac{\partial g}{\partial x_k} \frac{\partial x_i}{\partial x_i}$$

$$= A \frac{\partial x_j}{\partial x_i} + B \frac{\partial x_k}{\partial x_i} + C \frac{\partial x_i}{\partial x_i}$$

Finally let's apply our back prop rules to compute the gradient of our example expression.





Let's check our reverse-mode AD results with symbolic differentiation:

value of y at x=3, w=13: -9/10

partial derivatives at x=3, w=13: (-69/100)(9/100)

```
import sympy as sp
x, w = sp.symbols('x, w')
y = x^* 2/(x-w)
dydx = sp.diff(y,x)
dydw = sp.diff(y,w)
print('partial derivatives of y wrt x,w :')
                                                                                   xk = sin(xi)
display(dydx)
display(dydw)
values = \{x:3, w:13\}
                                                                                   dy/dxi =
print(f'value of y at x={values[x]}, w={values[w]}: {y.subs(values)}\n')
print(f'partial derivatives at x={values[x]}, w={values[w]}:', \
                                                                                     dy/dk * dxk/dxi
      f'{dydx.subs(values)}, {dydw.subs(values)}')
                                                                                    = "A" * cos(xi)
   partial derivatives of y wrt x,w :
```

# Quadrature

= numerical evaluation of definite integrals



Start with 1D:  $f:[a,b] \rightarrow \mathbb{R}$ . To approximate

The word quadrature comes from real-world estimating of

e.g. | | A = area of region R

Crude estimate: count the "quadrats".

area of stone = 6 to 9

$$J(f) = \int_{a}^{b} f = \int_{a}^{b} f \cos dx,$$

Why not just find a formula F such that F' = f, and use J(f) = F(b) - F(a)?

We will when we can, but a closed formula for such an F frequently (i) does not exist, or (ii) is a pain to obtain.

We will sample f at m+1 points  $\times_{\sigma}, \times_{\iota}, \cdots, \times_{m}$ , and devise a "quadrature rule"

$$Q(f) = q(f(x_0), f(x_1), ..., f(x_m))$$

 $\approx \mathcal{J}(f)$  as accurately as possible.

We would really like Q(cf) = c Q(f) for any constant c and Q(f+g) = Q(f)+Q(g)

to mirror the corresponding properties required of an integral.

This forces q to be a linear combination of its arguments (proof?):

$$Q(f) = \sum_{i=0}^{n} w_i f(x_i)$$

$$(b-a)$$
  $\underset{j=0}{\overset{m}{\sum}} \propto_j f(x_j)$  weights

 $Q(f) = \sum_{j=0}^{n} W_j f(x_j) \stackrel{\text{Convenience}}{=} (b-a) \sum_{j=0}^{n} x_j f(x_j)$  Factoring out (b-a) from the w<sub>j</sub>'s makes the  $x_j$ 's independent of the interval length.

It remains to choose the  $\{x_j\}$  and the  $\{x_j\}$  to get an accurate approximation of f

Based on our experience with interpolation, it seems likely that some placements of the  $x_j$  will be better than others. But supposed we've decided on the  $\{x_j\}$ . Maybe uniformly space, maybe not.

Then we need to decide on the weights

[wj] or equivalenty {a;}.

What principle(s) should we use to constrain the weights?

Example: m=0, one sample of f:

$$\{\times_j\} = (\times_o)$$

$$Q(f) = \omega_o f(x_o) = (b-a) \propto_o f(x_o).$$

What's a good choice for 💢 o 🕻

Seems like it is highly desirable to get the integral of a constant function exactly correct:

$$Q(1) = (b-a) \propto_0.$$

$$Q(1) = (b-a) \propto_{o} \cdot | \qquad \text{want} \\ = J(1) = \int_{a}^{b} dx = (b-a)$$

This regnires 00=1.

So our one-point quadrature rule is:

What if we have more than 1 sample point?

What does  $Q(1) = \int I dictate?$ 

$$Q(1) = (b-a) \sum_{j=0}^{m} \alpha_j f(x_j) \stackrel{f:1}{=} (b-a) \sum_{j=0}^{m} \alpha_j \cdot 1 \quad \text{med} \quad (b-a) = \int_a^b dx$$
That is, we need 
$$\sum_{j=0}^{m} \alpha_j \cdot \frac{1}{j} = 1 \quad \text{That is 1 constraint on m+1 variables.}$$

How about requiring the rule to integrate linear functions correctly?

$$Q(L) = \int L$$

That provides one more (linear) constraint on the weights.

Continuing the theme, we could require that Q gets the exact integral for all polynomials of degree m or less:

ξα<sub>i</sub>ζ, This is m+1 (linear) constraints on the m+1 variables

Such a Q is said to have "polynomial degree m".

We can arrive at the same quadrature rule Q in another way ...

P(f) e Pm be the unique polynomial of degree at most m

that interpolates the gata

 $\{(x_i,f(x_i))\}$ Then define Q(f) = p(f). (This integral is easy to do symbolically because p is a polynomial.)

We could write p(f) in terms of the Lagrange polynomials:
$$\bigcirc (f) = \int_{a}^{b} f(x)(x) dx = \int_{a}^{b} \int_{a}^{b} f(x) f(x) dx = \int_{a}^{b} \int_{a}^{b} f(x) dx =$$

We can show that the two definitions of Q are equivalent.

which can be computed symbolically exactly.

Another example:

(m+1)-point rule of polynomial degree m with

 $\times_{s}$  =  $\alpha$ ,  $\times_{m}$  =  $\beta$  and the others uniformly spaced in between.

This is called the "closed Newton-Cotes rule" of degree m.

Again note that once the  $\{x_i\}$  are chosen, the requirement to be of polynomial degree m fixes the weights.

If the m+1 nodes are uniformly spaced with the endpoints a,b omitted, this is called the "open Newton-Cotes" rule of degree m.

## Deriving the weights exactly using sympy

# derive weights for quadrature rule by integrating Lagrange form of interpolating polynomial import numpy as np

$$\begin{array}{l} xx = \text{sp.symbols}(x') \\ a,b = 0,2 \\ H = b-a \\ m = 3 \\ x = [\text{sp.Rational}(j^*H,m) \text{ for } j \text{ in } \text{range}(m+1)] \\ \text{display}(x) \\ \\ \text{e.c.} : & -\frac{9x^3}{16} + \frac{9x^2}{4} - \frac{11x}{4} + 1 \\ \text{o.c.} : & -\frac{9x^3}{16} - \frac{45x^2}{8} + \frac{9x}{2} \\ \text{logrange polymer for } i \text{ in } \text{range}(m+1): \\ \text{l.j.} = 1 \\ \text{for } j \text{ in } \text{range}(m+1): \\ \text{l.j.} = 1 \\ \text{for } i \text{ in } \text{range}(m+1): \\ \text{l.j.} = 1 \\ \text{l.j.} = \text{l.j.} = (xx-x[i])/(x[j]-x[i]) \\ \text{l.j.} = \text{l.j.} = (xx-x[i])/(x[j]-x[i]) \\ \text{l.j.} = \text{l.j.} = (xy-x[i])/(x[j]-x[i]) \\ \text{l.j.} = (xy-x[i])/(x[i]-x[i]) \\ \text{l.j.} = (xy-x[i])/(x[i]$$

$$\begin{cases} 7 \\ 3 \\ 3 \end{cases} = \begin{bmatrix} 0, 2/3, 4/3, 2 \end{bmatrix}$$

$$\begin{cases} l_0(x) : -\frac{9x^3}{16} + \frac{9x^2}{4} - \frac{11x}{4} + 1 \\ l_1(x) : \frac{27x^3}{16} - \frac{45x^2}{8} + \frac{9x}{2} \\ \\ l_2(x) = -\frac{27x^3}{16} + \frac{9x^2}{2} - \frac{9x}{4} \\ \\ l_3(x) = \frac{9x^3}{16} - \frac{9x^2}{8} + \frac{x}{2} \\ \\ 3(x) \\ \end{cases} = \begin{bmatrix} 1/8, 3/8, 3/8, 1/8 \end{bmatrix}$$

$$\begin{cases} 7 \\ 3 \\ 4 \end{cases}$$

Let's "funcify" the above code and look at some other concrete examples ...

# closed Newton-Cotes on several intervals open Newton-Cotes

#### randomly chosen points

#### Copyable code