

## Why derivatives?

If we want to minimize functions, we need to think about derivatives.

If  $x$  minimizes  $f(x)$ , then  $\frac{\partial f}{\partial x(i)} = 0$  for all  $i$

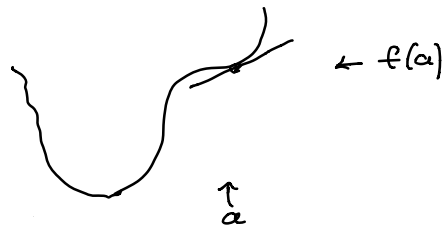
(necessary, but not sufficient)

If this does not hold, then can decrease  $f$ .

---

Recall, derivative of  $f$  at  $a$  gives a linear approximation of  $f$ .

In one variable



$$f(a+\delta) \approx f(a) + \delta f'(a)$$

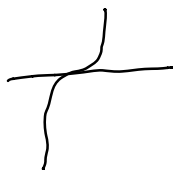
↑ linear in  $\delta$ .

Means  $\lim_{\delta \rightarrow 0} \frac{f(a+\delta) - f(a)}{\delta} \rightarrow f'(a)$

If  $f'(a) > 0$ , small negative  $\delta$  decreases  $f$   
 $< 0$ , small positive " " "

$f'(a)=0$  does not imply  $a$  is a minimum

$$f(x) = x^3$$



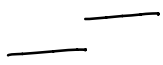
$$f'(0) = 0$$

Def

$f$  is continuous if  $f(x) \rightarrow f(y)$   
as  $x \rightarrow y$

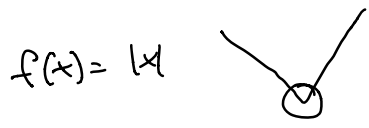
not continuous

$$f(x) = \begin{cases} 1 & x > 0 \\ 0 & x \leq 0 \end{cases}$$



$f$  is differentiable if  $f'(x)$   
exists for all  $x$

not differentiable at 0



$f$  is continuously differentiable  
if  $f'(x)$  is a continuous function

skip

$$f(x) = \begin{cases} x^2 \sin(1/x) & x \neq 0 \\ 0 & x = 0 \end{cases}$$

strange if continuous  
& differentiable but not  $C^1$

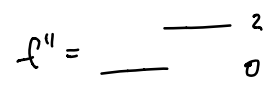
$f$  is in  $C^k$  if first  $k$   
derivatives exist and are continuous.

$C^1$  but not  $C^2$

$$f(x) = \begin{cases} x^2 & x > 0 \\ 0 & x = 0 \end{cases}$$



Branches and comparizers break continuity.  
But only at a few points.



we care about derivatives where they exist

In many variables,  $x \in \mathbb{R}^n$ ,  $\delta \in \mathbb{R}^n$

Vary  $x(i)$  while treating other variables as constants

$$f(a(i) + \delta(i), a(2), \dots, a(n)) - f(a(1), \dots, a(n)) \\ \approx \delta(i) \frac{\partial f}{\partial x(i)}(a)$$

$$f(a + \delta) \approx f(a) + \sum_i \delta(i) \frac{\partial f}{\partial a(i)} \\ = f(a) + (\nabla f)(a)^T \delta$$

$(\nabla f)(a)$  is the column vector  $\left( \frac{\partial f}{\partial x(1)}(a), \dots, \frac{\partial f}{\partial x(n)}(a) \right)$

Called the gradient.

$f: \mathbb{R}^n \rightarrow \mathbb{R}$  is differentiable at  $a$  if  $\exists g \in \mathbb{R}^n$  st.

$$\lim_{\delta \rightarrow 0} \sup_{\|\delta\|_2 \leq \epsilon} \frac{|f(a + \delta) - (f(a) + g^T \delta)|}{\|\delta\|} \rightarrow 0$$

In which case  $g = (\nabla f)(a)$

That is, the linear approximation is good,  
and improves for smaller  $\delta$ .

---

Automatic differentiation: Code that computes derivatives  
Not symbolic, not a numerical approximation

## Forward mode

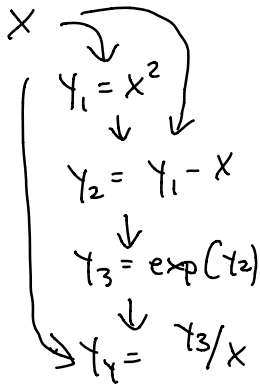
Compute intermediate terms in order,  $y_1, \dots, y_m$

$y_i = f_i(x_1, \dots, x_n, y_1, \dots, y_{i-1})$  depends on previous.

Usually,  $f_i$  depends on one or two of these.

Eg.  $f(x) = \frac{\exp(x^2 - x)}{x}$

compute  $\frac{dy_i}{dx}$  in order,  
using previous.



$$\frac{dy_1}{dx} = 2x$$

$$\frac{dy_2}{dx} = \frac{dy_1}{dx} - \frac{dx}{dx} = 2x - 1$$

$$\frac{dy_3}{dx} = \frac{d \exp(t_2)}{dx} = \exp(t_2) \cdot \frac{dy_2}{dx}$$

and so on.

Can turn (almost) any program into this form,

for any input.

Will get different circuits/graphs for different inputs

because of branches.

Continuous when branch conditions are not tight.

# of steps here is # of operations

Example

func(x, k)

$$z = 1$$

for i in 1 to k

$$z = x + \frac{1}{z}$$

return z

Only differentiable in x.

If call with func(x, 3) get

$$y_1 = 1$$

$$y_2 = \frac{1}{y_1}$$

$$y_3 = x + y_2$$

$$y_4 = \frac{1}{y_3}$$

$$y_5 = x + y_4$$

$$y_6 = \frac{1}{y_5}$$

$$y_7 = x + y_6$$

key idea: for each function  $\gamma_i$

only need to compute  $\gamma_i(a)$  and  $\gamma_i'(a) = \frac{d}{dx} \gamma_i(a)$

Examples: If

$$\gamma_k = \gamma_i + \gamma_j \quad \gamma_k' = \gamma_i' + \gamma_j'$$

$$\gamma_k = \gamma_i \cdot \gamma_j \quad \gamma_k' = \gamma_i \gamma_j' + \gamma_j \gamma_i'$$

$$\gamma_k = \frac{\gamma_i}{\gamma_j} \quad \gamma_k' = \frac{\gamma_i \gamma_j' - \gamma_i' \gamma_j}{\gamma_j^2}$$

$$\gamma_k = \exp(\gamma_i) \quad \gamma_k' = \exp(\gamma_i) \gamma_i'$$

Chain rule says  $(f(g(x)))' = f'(g(x)) \cdot g'(x)$

$$\text{so } \gamma_k = g(\gamma_i) \quad \gamma_k' = g'(\gamma_i) \cdot \gamma_i'$$

In fact, can handle a multi-input  $g$

$$\gamma_k = g(\gamma_i, \dots, \gamma_j)$$

By multivariate chain rule

$$\begin{aligned} \frac{d}{dx} \gamma_k(x) &= \frac{d}{dx} g(\gamma_i(x), \dots, \gamma_j(x)) \\ &= \sum_{h=i}^j \frac{\partial g}{\partial \gamma_h} \frac{d\gamma_h}{dx} \end{aligned}$$

Example:

$$Y_1 = X$$

$$Y_2 = X$$

$$Y_3 = Y_2$$

$$Y_4 = Y_1 \cdot Y_2 \cdot Y_3$$

$$\frac{dY_4}{dx} = \frac{\partial Y_4}{\partial Y_1} \cdot \frac{dY_1}{dx} + \frac{\partial Y_4}{\partial Y_2} \cdot \frac{dY_2}{dx} + \frac{\partial Y_4}{\partial Y_3} \cdot \frac{dY_3}{dx}$$

$$= Y_2 Y_3 \cdot \frac{dY_1}{dx} + Y_1 Y_3 \frac{dY_2}{dx} + Y_1 Y_2 \frac{dY_3}{dx}$$

"  
1

"  
1

"  
1

$$\frac{dY_3}{dY_2} \cdot \frac{dY_2}{dx} = 1$$

$$= Y_2 Y_3 + Y_1 Y_3 + Y_1 Y_2$$

$$= 3x^2$$

Let's see why is true.

First, in one variable

$$\begin{aligned} f(\gamma(a+\delta)) &\approx f(\gamma(a) + \delta \gamma'(a)) \\ &\approx f(\gamma(a)) + f'(\gamma(a)) \cdot \delta \gamma'(a) \\ &= f(\gamma(a)) + \delta \cdot \underline{f'(\gamma(a)) \cdot \gamma'(a)} \end{aligned}$$

In many

$$\begin{aligned} f(\gamma_1(a+\delta), \dots, \gamma_n(a+\delta)) \\ &\approx f(\gamma_1(a) + \delta \gamma_1'(a), \dots, \gamma_n(a) + \delta \gamma_n'(a)) \\ &\approx f(\gamma_1(a), \dots, \gamma_n(a)) + \underbrace{\sum_{i=1}^n \frac{\partial f}{\partial \gamma_i} \cdot \delta \gamma_i'(a)} \end{aligned}$$

If  $\gamma: \mathbb{R} \rightarrow \mathbb{R}^n$   $a \rightarrow (\gamma_1(a), \dots, \gamma_n(a))$

$$\gamma'(a) = (\gamma_1'(a), \dots, \gamma_n'(a))$$

can write as  $(\nabla f)(a)^T \gamma$

Forward mode: compute  $f(a)$  and  $f'(a)$

in essentially 3 times as many operations as to compute  $f(a)$  alone.



For  $f: \mathbb{R}^n \rightarrow \mathbb{R}$ ,  $f(x_1, \dots, x_n)$

compute  $f(x)$  and  $\frac{\partial f}{\partial x_1}$

then  $f(x)$  and  $\frac{\partial f}{\partial x_2}$

etc.

Get  $f(x)$ ,  $(\nabla f)(x)$  in  $3n$  times as many steps as for  $f(x)$ .

Can reduce cost by computing all simultaneously.

But uses more memory.

---

Reverse mode saves a factor of  $n$  !

Reverse mode uses two passes.

forward: compute  $y_i(a)$  for every  $i$ , as usual.

and  $\frac{\partial y_j}{\partial x_i}(a)$  for every  $i$  that is an input to  $j$   
(usually 1 or 2)

Then, a reverse pass to compute

$\frac{\partial f}{\partial y_i}$  for each  $i$ , from  $n$  down to 1

and then  $\frac{\partial f}{\partial x_i}$

What does this mean?

$$f = \gamma_m, \text{ so } \frac{\partial f}{\partial \gamma_m} = 1$$

When I say  $\frac{\partial f}{\partial \gamma_i}$  I mean to view

$f$  as a function of  $x_1, \dots, x_n$  and  $\gamma_1, \dots, \gamma_i$   
treating all of those as variables.

And, for  $j > i$   $\gamma_j = f_j(x_1, \dots, x_n, \gamma_1, \dots, \gamma_i)$  (\*)

Formally, let  $F_i(x_1, \dots, x_n, \gamma_1, \dots, \gamma_i)$

be the function determined by (\*), output is  $\gamma_m = f$ .

We compute  $\frac{\partial F_i}{\partial \gamma_i}$ , and finish with  $\frac{\partial F_0}{\partial x_i}$  for each  $i$

To compute  $\frac{\partial F_i}{\partial \gamma_i}$ , treat  $x_1, \dots, x_n, \gamma_1, \dots, \gamma_{i-1}$  as fixed,

so  $F_i = h(f_k(\bar{x}, \gamma_1, \dots, \gamma_i))$  for  $f_k$  taking  $i$  as input

The chain rule tells us that

$$\frac{\partial F_i}{\partial \gamma_i} = \sum_{k \text{ st. } i \text{ is input to } k} \frac{\partial F}{\partial \gamma_k} \frac{\partial \gamma_k}{\partial \gamma_i}$$

$\frac{\partial \gamma_k}{\partial \gamma_i}$  was computed in forward pass

$k > i$ , so  $\frac{\partial F}{\partial \gamma_k}$  computed earlier in backward pass.

Example  $f(x_1, x_2) = x_1 \cdot \exp(x_2) - x_1$

$$\gamma_1 = \exp(x_2)$$

$$\gamma_2 = x_1 \cdot \gamma_1$$

$$\gamma_3 = \gamma_2 - x_1$$

forwards, with input  $a_1, a_2$

compute  $\gamma_1 = \exp(a_2) \quad \frac{d\gamma_1}{dx_2} = \exp(a_2)$

$$\gamma_2 = a_1 \cdot \exp(a_2) \quad \frac{d\gamma_2}{d\gamma_1} = a_1 \quad \frac{d\gamma_2}{dx_1} = \gamma_1$$
$$= a_1 \quad = \exp(a_2)$$

$$\gamma_3 = a_1 \cdot \exp(a_2) - a_1 \quad \frac{d\gamma_3}{d\gamma_2} = 1 \quad \frac{d\gamma_3}{dx_1} = -1$$

Backwards:  $\frac{\partial f}{\partial \gamma_3} = 1$

$$\frac{\partial f}{\partial \gamma_2} = \frac{\partial f}{\partial \gamma_3} \cdot \frac{\partial \gamma_3}{\partial \gamma_2} = 1 \cdot 1 = 1$$

$$\frac{\partial f}{\partial \gamma_1} = \frac{\partial f}{\partial \gamma_2} \cdot \frac{\partial \gamma_2}{\partial \gamma_1} = 1 \cdot a_1 = a_1$$

$$\frac{\partial f}{\partial x_2} = \frac{\partial f}{\partial \gamma_1} \cdot \frac{\partial \gamma_1}{\partial x_2} = a_1 \cdot \exp(a_2)$$

$$\frac{\partial f}{\partial x_1} = \frac{\partial f}{\partial \gamma_2} \cdot \frac{\partial \gamma_2}{\partial x_1} + \frac{\partial f}{\partial \gamma_3} \cdot \frac{\partial \gamma_3}{\partial x_1} = 1 \cdot \gamma_1 + 1 \cdot (-1) = \gamma_1 - 1$$
$$\exp(a_2) - 1$$

Problems: have to store everything,  
so memory use is proportional to time.  
Are ways to make that better.

Stopped Here

Computing derivatives in one variable.

Could use definition

$$f'(a) = \lim_{\delta \rightarrow 0} \frac{f(a+\delta) - f(a)}{\delta} \quad \text{for some } \delta.$$

Precision issues, let  $\tilde{f}$  be what we compute in floating point, and  $u$  be precision.

$$\text{So, } |f(a) - \tilde{f}(a)| \approx u |f(a)|$$

$$|f(a+\delta) - \tilde{f}(a+\delta)| \approx u |f(a+\delta)| \approx u |f(a)|$$

Let  $L$  be st.

$$f(a+\delta) \approx f(a) + \delta f'(a) + \delta^2 L \quad (\text{Taylor series})$$

$$\text{So, } \left| \frac{\tilde{f}(a+\delta) - \tilde{f}(a)}{\delta} - f'(a) \right|$$

$$\approx \left| \frac{f(a+\delta) - f(a)}{\delta} - f'(a) \right| + \frac{2u |f(a)|}{\delta}$$

$$\approx \frac{L \delta^2}{\delta} + \frac{2u |f(a)|}{\delta} \approx L \delta + \frac{2u |f(a)|}{\delta}$$

$$\text{minimized when } \delta^2 = \frac{u |f(a)|}{L} \quad \delta^2 \sim \sqrt{u}$$