· loss.backward() $\longrightarrow$ param.grad $\leftarrow$ what grad? why?

$1.bw$?

# Computational Differentiation

① Numerical Differentiation

$f(\underline{x})$  $\frac{\partial f(\underline{x})}{\partial \underline{x}} = \left[ \frac{\partial f(\underline{x})}{\partial x_1}, \cdots \frac{\partial f(\underline{x})}{\partial x_n} \right]$

<span style="color:red">very good test case to verify the correctness of your implementation</span>

Disadvantages:

a) Approximate

$\frac{\partial f(\underline{x})}{\partial x_i} = \frac{f(\underline{x} + \underline{e}_i \varepsilon) - f(\underline{x})}{\varepsilon}$   small no $\approx 10^{-4}, 10^{-6}$

b) Computationally expensive:

$\underline{e}_i = \begin{bmatrix} 0 \\ \vdots \\ i \\ \vdots \\ 0 \end{bmatrix} \leftarrow$ 1 at $i^{th}$ place, 0 everywhere

How many times do you have to call

$f(\underline{x})$ to compute $\frac{\partial f(x)}{\partial \underline{x}}$ if $\underline{x} \in \mathbb{R}^n$?

$n+1$ times

② Symbolic Differentiation (SD)

e.g. sympy / Matlab

Disadvantages

a) you need a new language to communicate your formulass

b) You get formula as output. Efficiency is not considered

formula / code $\longrightarrow$ code to implement the derivative comp. complex?

formula of a function $\longrightarrow$ [SD] $\longrightarrow$ formula for the derivative

③ Automatic differentiation (AD)

→ Operator overloading to implement (AD)

C++ / Python

$a + b$

$\underset{int}{\uparrow} \quad \underset{int}{\uparrow}$

$a = Tensor()$
$b = Tensor()$

$a = Person()$
$b = Person()$

$a + b$

$a @ b$

Person.__add__

Person.__matmul__

```
class Tensor( )
    def __init__( )
        self.value =
        self.grad =

a = tensor(int)
b = Tensor(int)
```

$c = a + b$

$c.value = a.value + b.value$

$c.grad = \underline{\hspace{3cm}}$

```
a = Sympy.Variable ("a")
b =       "      "      ("b")
```

$c = a + b$

$c.expression = $ "a+b"
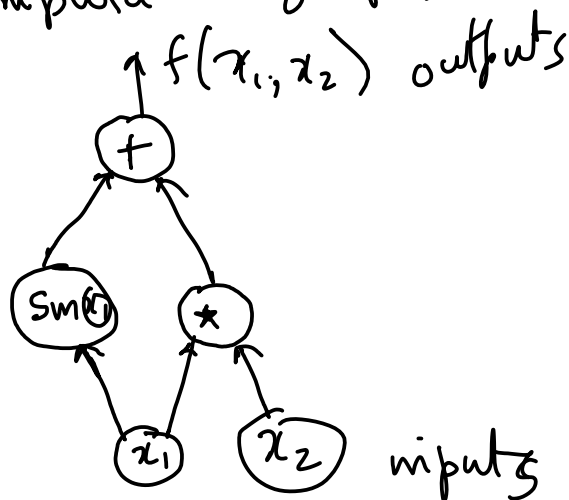
---

AD

① Forward-mode AD (Forward accumulative)

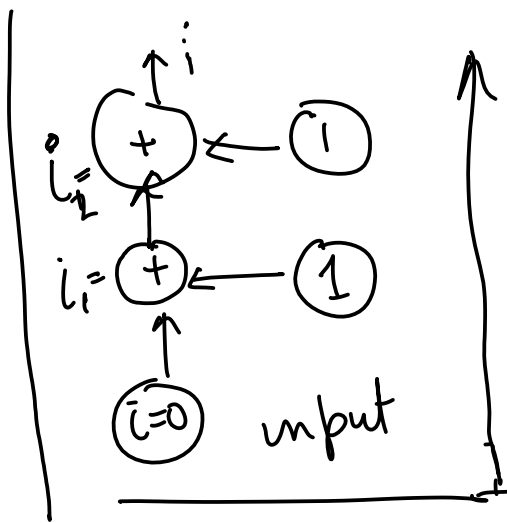② Reverse-mode AD (Backpropagation)
                         in NN

Mathematical expressions can be represented as
a <u>Directed Acyclic graph</u> (computation graph)

Example

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

$i = \text{Tensor}(0)$
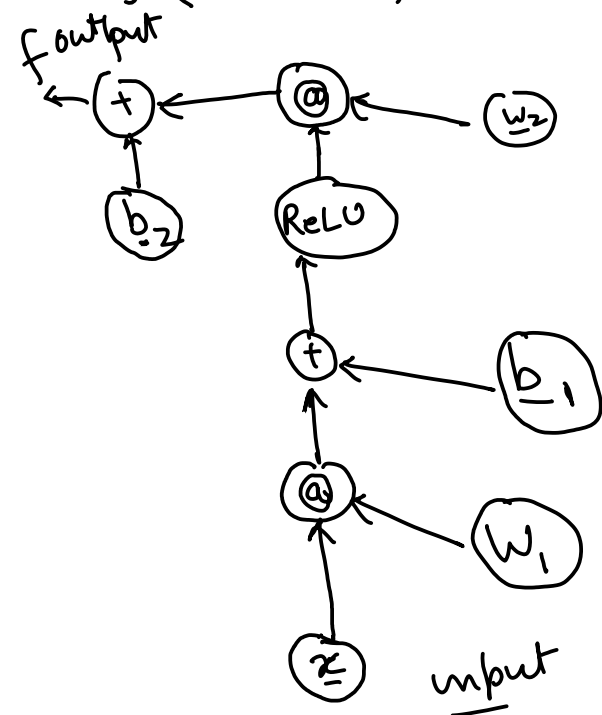for $i$ in range(10)
$\quad i = i + 1$



Internally
Pytorch
will
keep a
track of
this graph

## Computation Graph of MLP (2-layer, ReLU)

$$f(\underline{x}; \underline{W}_1, \underline{w}_2, \underline{b}_1, \underline{b}_2) = \underline{w}_2^T \text{ReLU}(W_1 \underline{x} + \underline{b}_1) + b_2$$

$$= \underline{w}_2^T (W_1 \underline{x} + b_1)_+ + \underline{b}_2$$

ReLU

# ① Forward-mode AD
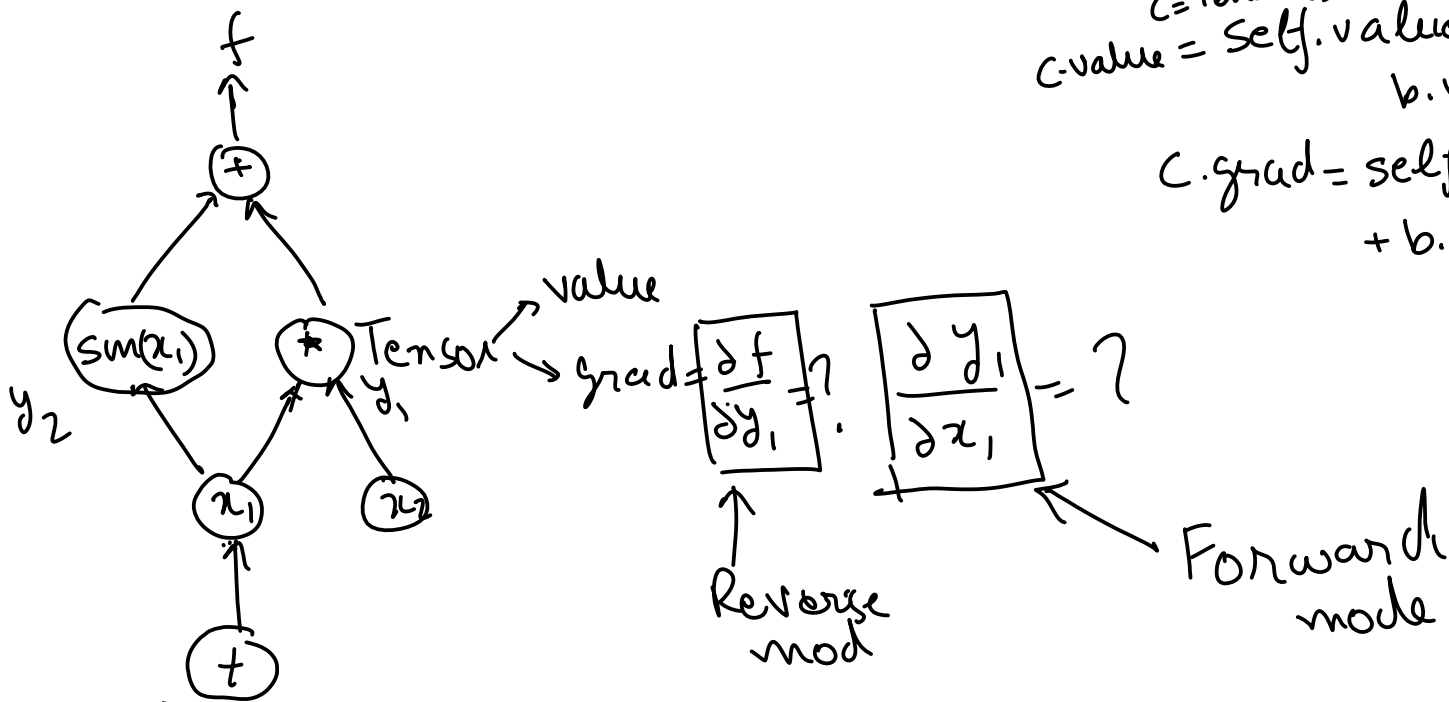
Use operator overloading to implement AD

```
class Tensor()
    def __init__( )
        self.value =
        self.grad =
```

$a = tensor(int)$
$b = Tensor(int)$

$c = a + b$



class Tensor
```
def __add__(s, b)
    (c = Tensor()
    c.value = self.value +
                    b.value
    c.grad = self.grad
                + b.grad
```

value
grad $= \dfrac{\partial f}{\partial y_1}$ ?  $\boxed{\dfrac{\partial y_1}{\partial x_1}} = ?$

Reverse mod ↑

Forward mode

In forward mode   $x_1.grad = \dfrac{\partial x_1}{\partial t} = 1$

$x_2.grad = \dfrac{\partial x_2}{\partial t} = 0$

$y_1 = x_1 * x_2 \Rightarrow y_1.grad = \dfrac{\partial(x_1^* x_2)}{\partial t} = x_1 \dfrac{\partial x_2}{\partial t} + x_2 \dfrac{\partial x_1}{\partial t}$

$$\Rightarrow y_1.grad = x_1*(x_2.grad) + x_2*(x_1.grad)$$

$$y_2 = \sin(x_1) \quad \rightarrow \frac{\partial y_2}{\partial t} = y_2.grad = \cos(x_1) \frac{\partial x}{\partial t} = \cos(x_1) [x_1.grad]$$
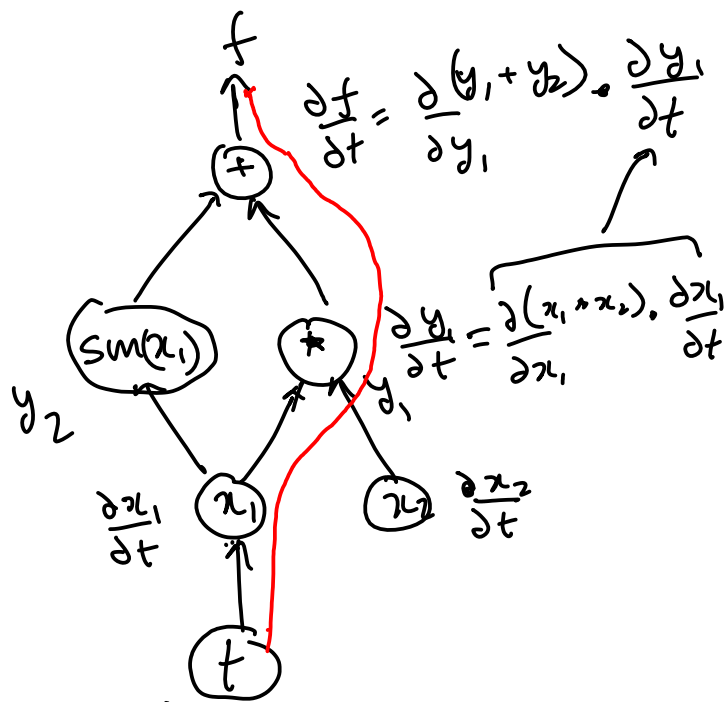
$$f = y_1 + y_2 \qquad \Rightarrow \frac{\partial f}{\partial t} = f.grad = \frac{\partial y_1}{\partial t} + \frac{\partial y_2}{\partial t}$$

$$= y_1.grad + y_2.grad$$

---

## Accumulation direction

$$\frac{\partial f}{\partial t} = \left[ \frac{\partial}{\partial y_1} (y_1 + y_2) . \left[ \frac{\partial (x_1 * x_2)}{\partial x_1} . \frac{\partial x_1}{\partial t} \right] \right]$$

This accumulation happens first



$$\frac{\partial f}{\partial t} = \frac{\partial (y_1 + y_2)}{\partial y_1} . \frac{\partial y_1}{\partial t}$$

$$\frac{\partial y_1}{\partial t} = \frac{\partial (x_1 * x_2)}{\partial x_1} . \frac{\partial x_1}{\partial t}$$

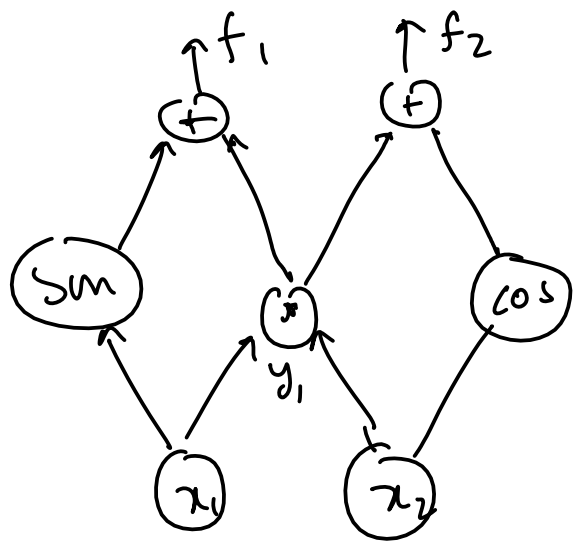$$f = h\Big(g\big(k(\ell(x))\big)\Big)$$

$$\frac{\partial f}{\partial x} = \left( \frac{\partial h}{\partial g} \left( \frac{\partial g}{\partial k} \left( \frac{\partial k}{\partial \ell} \frac{\partial \ell}{\partial x} \right) \right) \right)$$

① In forward mode differentiation, the accumulation of chain rule derivatives happens from the input side to the output side (forward direction).

② In forward mode differentiation, you have to re-propagate the derivatives for each input.
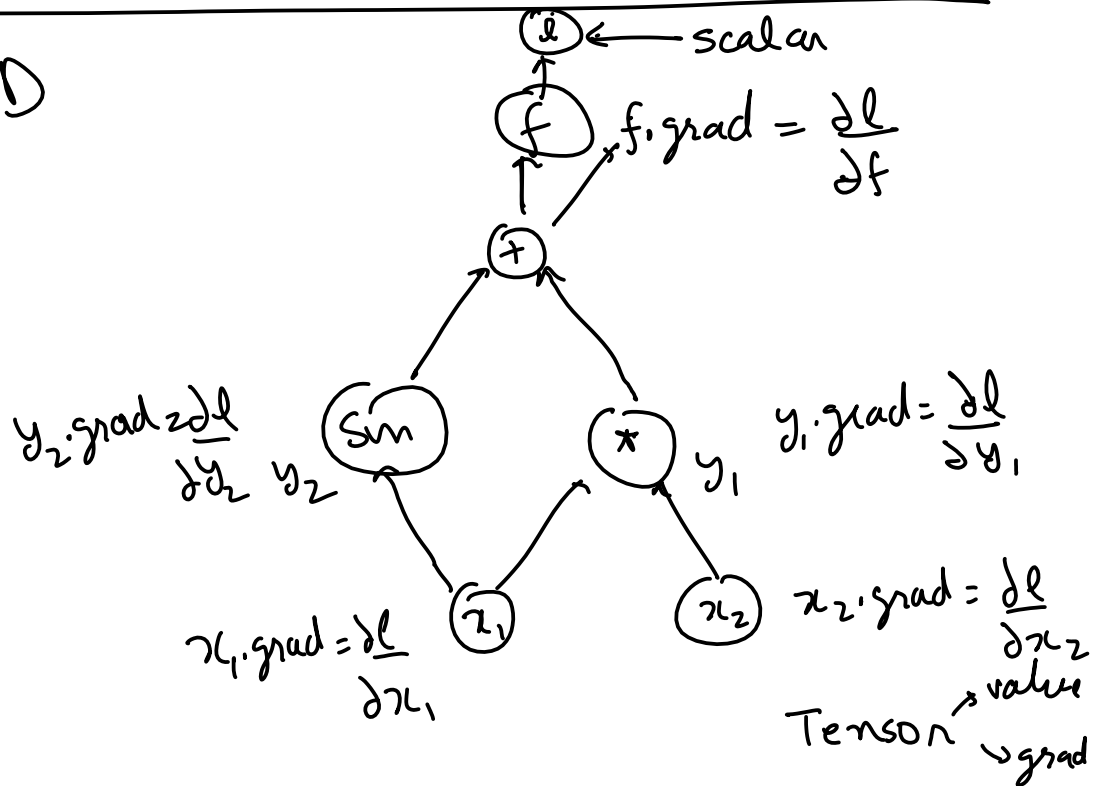
$$f_1(x_1, x_2) = x_1 x_2 + \sin(x_1)$$

$$f_2(x_1, x_2) = \underbrace{x_1 \cdot x_2} + \cos(x_2)$$

$$\frac{\partial f_1}{\partial t} \quad , \quad \frac{\partial f_2}{\partial t}$$

③ we dont have to recompute derivatives for the common part of the graph.

---

② Reverse - mode AD

$\ell \leftarrow$ — scalar

$f.grad = \dfrac{\partial \ell}{\partial f}$

$y_2.grad = \dfrac{\partial \ell}{\partial y_2} \quad y_2$

$y_1.grad = \dfrac{\partial \ell}{\partial y_1}$

$x_1.grad = \dfrac{\partial \ell}{\partial x_1}$

$x_2.grad = \dfrac{\partial \ell}{\partial x_2}$

Tensor $\nearrow$ value $\searrow$ grad

Before you turn this problem in, make sure everything runs as expected. First, **restart the kernel** (in the menubar, select Kernel→Restart) and then **run all cells** (in the menubar, select Cell→Run All).

Make sure you fill in any place that says `YOUR CODE HERE` or "YOUR ANSWER HERE", as well as your name and collaborators below:

```
In [ ]:  NAME = ""
         COLLABORATORS = ""
```
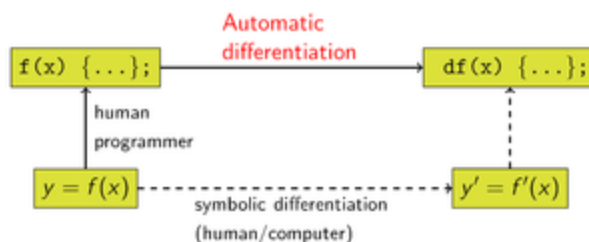
---

# Differentiation options

1. Numerical differentiation

2. Symbolic differentiation

3. Automatic differentiation

    A. Forward mode differentiation
    B. Reverse mode differentiation

# 1. Numerical differentiation

# 2. Symbolic differentiation

# 3. Automatic differentiation



## 3.A Forward mode

Example:

$$z = f(x_1, x_2) = x_1 x_2 + \sin(x_1)$$

## 3.B Reverse mode

Example:

$$z = f(x_1, x_2) = x_1 x_2 + \sin(x_1)$$

```python
In [ ]: import numpy as np
        class ForwardDiff:
            def __init__(self, value, grad=None):
                self.value = value
                self.grad = np.zeros_like(value) if grad is None else grad


            def __add__(self, other):
                cls = type(self)
                other = other if isinstance(other, cls) else cls(other)
                out = cls(self.value + other.value,
                          self.grad + other.grad)
                return out
            __radd__ = __add__

            def __repr__(self):
                return f"{self.__class__.__name__}(data={self.value}, grad={self.gra
        x = ForwardDiff(2, 1)
        y = ForwardDiff(3, 0)

        f = x + y
        f
```

```python
In [ ]: oldFD = ForwardDiff # Bad practice: do not do it
        class ForwardDiff(oldFD):
            def __mul__(self, other):
                cls = type(self)
                other = other if isinstance(other, cls) else cls(other)
                out = cls(self.value * other.value,
                          other.value * self.grad+
                          self.value * other.grad)
                return out

            __rmul__ = __mul__

        x = ForwardDiff(2, 0)
        y = ForwardDiff(3, 1)

        f1 = x * y
        f2 = 2*x + 3*y + x*y
        f1, f2
```

```python
In [ ]: oldFD = ForwardDiff # Bad practice: do not do it
        class ForwardDiff(oldFD):
            def log(self):
                cls = type(self)
                return cls(np.log(self.value),
                           1/self.value * self.grad)
            def exp(self):
                cls = type(self)
```

```python
        out_val = np.exp(self.value)
        return cls(out_val,
                   out_val * self.grad)

    def sin(self):
        cls = type(self)
        return cls(np.sin(self.value),
                   np.cos(self.value) * self.grad)

    def cos(self):
        cls = type(self)
        return cls(np.cos(self.value),
                   -np.sin(self.value) * self.grad)

    def __pow__(self, other):
        cls = type(self)
        other = other if isinstance(other, cls) else cls(other)
        return (self.log() * other).exp()

    def __neg__(self): # -self
        return self * -1

    def __sub__(self, other): # self - other
        return self + (-other)

    def __truediv__(self, other): # self / other
        return self * other**-1

    def __rtruediv__(self, other): # other / self
        return other * self**-1


x = ForwardDiff(2, 1)
y = ForwardDiff(3, 0)

f = x**y
f
```

```python
import numpy as np
def add_vjp(a, b, grad):
    return grad, grad


def no_parents_vjp(grad):
    return (grad,)

class ReverseDiff:
    def __init__(self, value, parents=(), op='', vjp=no_parents_vjp):
        self.value = value
        self.parents = parents
        self.op = op
        self.vjp = vjp
        self.grad = None

    def backward(self, grad):
```

```
            self.grad = grad
            op_args = [p.value for p in self.parents]
            grads = self.vjp(*op_args, grad)
            for g, p in zip(grads, self.parents):
                p.backward(g)

    def __add__(self, other):
        cls = type(self)
        other = other if isinstance(other, cls) else cls(other)
        out = cls(self.value + other.value,
                  parents=(self, other),
                  op='+',
                  vjp=add_vjp)
        return out

    __radd__ = __add__

    def __repr__(self):
        cls = type(self)
        return f"{cls.__name__}(value={self.value}, parents={self.parents},

x = ReverseDiff(2)
y = ReverseDiff(3)

f = x + y + 3
f.backward(1)
f
x.grad, y.grad
```

In [ ]:
```
oldRD = ReverseDiff # Bad practice: do not do it

def mul_vjp(a, b, grad):
    return grad * b, grad * a

class ReverseDiff(oldRD):
    def __mul__(self, other):
        cls = type(self)
        other = other if isinstance(other, cls) else cls(other)
        out = cls(self.value * other.value,
                  parents=(self, other),
                  op='*',
                  vjp=mul_vjp)
        return out

    __rmul__ = __mul__

x = ReverseDiff(2)
y = ReverseDiff(3)

f1 = 5*x + 7* y
f1.backward(1)
x.grad, y.grad
```

In [ ]:
```
f2 = x*y
f2.backward(1)
```

```
x.grad, y.grad
```