

# AC 207 Final Project: Milestone 2

Group 22: Neil Sehgal, Lotus Xia, Andrew Zhang, Diwei Zhang

November 2021

## 1 Introduction

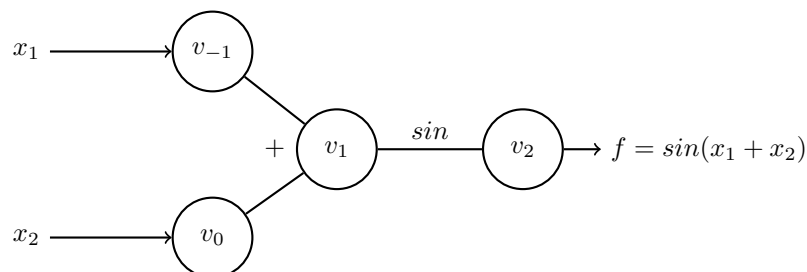
Differentiation is a central operation in science and engineering, used in various settings from financial markets to physics labs. For example, to optimize a function we must find the minima or maxima by using differentiation tests. There are multiple computational mechanisms to find derivatives. For example, numerical approximations like the Finite-Differences method involve utilizing the definition of a derivative and plugging in small values of  $h$  to evaluate the function. However, this method suffers from poor accuracy and stability issues due to floating point errors. Symbolic computation is another method, more accurate than numerical approximations, but often not applicable due to its unweildly and complex overhead. Here, we focus on automatic differentiation (AD), which overcomes the drawbacks of these two issues: it has less overhead than symbolic differentiation while calculating derivatives at machine precision. This makes it an integral part of machine learning algorithms, which demand computational efficiency as well as high precision.

## 2 Background

Automatic differentiation can be approached in two ways, known as *forward mode* and *reverse mode*. Both methods depend on the decomposition of a complex function into elementary operations. Then, we can take advantage of the chain rule to compute the derivatives of complex functions from trivially computed derivatives of these elementary functions.

### 2.1 Computational Graph

The first step of AD is to construct a computational graph which relates the elementary inputs to the final function. Each node represents an intermediate computation  $v_n$ , and edges connect nodes via elementary functions  $V_n$ . For example, given a multivariate scalar function  $f(x_1, x_2) = \sin(x_1 + x_2)$ , the computational graph looks like the following:



Elementary functions include but are not limited to the following: algebraic operations, exponentials, trigonometrics, and their inverses. AD takes advantage of the fact that derivatives of elementary functions are trivial to compute, allowing the differentiation of complex functions via the chain rule.

## 2.2 Chain Rule

Suppose we want to calculate  $\frac{\partial f}{\partial x_1}$ . From the computational graph, we can rewrite  $f$  in terms of nested intermediate functions:

$$f(x_1, x_2) = V_2(V_1(V_{-1}(x_1), V_0(x_2)))$$

where  $V_n$  is the function that gives the value of node  $v_n$ . The chain rule states that

$$\frac{\partial f}{\partial x_1} = \frac{\partial V_2}{\partial V_1} \frac{\partial V_1}{\partial V_{-1}} \frac{\partial V_{-1}}{\partial x_1}$$

For a multivariate function  $f(g_1(x), g_2(x), \dots, g_n(x))$ , the chain rule generalizes to:

$$\frac{\partial f}{\partial x} = \sum_{i=1}^n \left( \frac{\partial f}{\partial g_i} \frac{\partial g_i}{\partial x} \right)$$

## 2.3 Forward Mode

Forward mode, which is the conceptually simpler method, involves computing the chain rule from inside to outside. For the above example  $f(x_1, x_2) = \sin(x_1 + x_2)$ , the *evaluation trace* at the point  $(1, 2)$  is generated using the computational graph as a guide:

trace	elem	val	elem_der	$\nabla_{x_1}$	$\nabla_{x_2}$
$v_{-1} = x_1$	$x_1$	1	$v_{-1}$	1	0
$v_0 = x_2$	$x_2$	2	$v_0$	0	1
$v_1$	$v_{-1} + v_0$	3	$v_{-1} + v_0$	1	1
$v_2 = f(x_1, x_2)$	$\sin(v_1)$	$\sin(3)$	$\cos(v_1)v_1$	$\cos(3)$	$\cos(3)$

The evaluation trace essentially computes a value (primal trace) and derivative (tangent trace) for each node of the computational graph. The desired final results are  $\frac{\partial f}{\partial x_1}$  and  $\frac{\partial f}{\partial x_2}$ , shown in the last row of the trace using the gradient symbol  $\nabla$ . Forward mode is advantageous when there are few inputs ( $m$ ) and many outputs ( $n$ ). The opposite is true for reverse mode, which motivates its use in deep learning applications with many inputs.

## 2.4 Gradient & Jacobian

The gradient of a multivariate scalar-valued function  $f(g_1, g_2, \dots, g_m)$  gives all of its partial derivatives in a vector of size  $m \times 1$ :

$$\nabla f = \nabla f_g = \begin{bmatrix} \frac{\partial f}{\partial g_1} \\ \frac{\partial f}{\partial g_2} \\ \vdots \\ \frac{\partial f}{\partial g_m} \end{bmatrix}$$

If the inputs  $g_i$  are themselves multivariate functions with inputs  $x \in \mathbb{R}^m$ , the chain rule gives

$$\nabla_x f = \sum_{i=1}^n \left( \frac{\partial f}{\partial g_i} \nabla g_i(x) \right)$$

Finally, we reach the most general form of the gradient: for multivariate vector-valued functions with  $m$  inputs and  $n$  outputs, the derivative of outputs with respect to inputs can be represented in a *Jacobian matrix* of dimensions  $n \times m$ :

$$f(x_1, x_2, \dots, x_m) = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

$$J_f = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \cdots & \frac{\partial y_1}{\partial x_m} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_n}{\partial x_1} & \cdots & \frac{\partial y_n}{\partial x_m} \end{bmatrix}$$

## 2.5 Dual Numbers

Dual numbers offer a convenient way to store and compute the primal and tangent traces in forward mode AD. A dual number has the form

$$a + b\epsilon,$$

where  $a, b \in \mathbb{R}$  and  $\epsilon^2 = 0$ . Combined with Taylor series expansion of a function  $f$  centered at  $a$ , dual numbers give rise to the property

$$f(a + b\epsilon) = f(a) + f'(a)b\epsilon$$

because all terms with powers of  $\epsilon$  higher than 2 are equal to zero. Thus, the Taylor series approximation of a function about a dual number gives both its value and derivative, in a single calculation. In forward mode AD it is convenient to represent both the primal and tangent trace of each node as a dual number, where the real part corresponds to the primal and the dual part corresponds to the tangent. As shown below, functions (operations) on these nodes preserve the relationship between the real and dual parts:

$$g(f(a + \epsilon)) = g(f(a) + f'(a)\epsilon) = g(f(a)) + g'(f(a))f'(a)\epsilon$$

Note that the coefficient of  $\epsilon$  is simply the derivative of the real part  $g(f(a))$ , as given by the chain rule.

## 3 How to Use Salad

### 3.1 Installation

The users will pull **Salad** from PyPI using command line

---

```
pip install salad
```

---

Then import **salad** and other dependencies (e.g. **numpy**) in python script

---

```
import salad as ad
import numpy as np
```

---

### 3.2 Methods

Users can use our package to deal with 3 different scenarios: 1 input to 1 output; multiple inputs to 1 output; multiple inputs to multiple outputs.

### 3.2.1 1 input to 1 output

Users may specify a single function with one input variable. See the example below with  $f(x) = \exp(2x)$  with  $x = 2$ .

---

```
x = ad.Variable(2, label="x") # specifies the input variable
f = ad.exp(2 * x) # specifies functions involving the variable
```

---

We make two observations at this stage: First, arithmetic operations (exp in our example) need to be implemented within `salad`. Second, `f` is of type `Variable`.

To retrieve the value and derivative, call

---

```
val = f.val # array(exp(4)), numpy.array
der = f.der["x"] # 2exp(4), float
```

---

### 3.2.2 m inputs to 1 output

Users may also specify a single function with many input variables. See the example below with  $f(x, y) = \exp(x + y)$ , with  $x = 1, y = 1$ .

---

```
x = ad.Variable(1, label='x') # specifies input variable x
y = ad.Variable(1, label='y') # specifies input variable y
f = ad.exp(x + 2 * y) # specifies functions involving the variables

val = f.val # array(exp(3)), numpy.array
dx = f.der["x"] # exp(3), float
dy = f.der["y"] # 2exp(3), float
grad = f.der # {"x": exp(3), "y": 2exp(3)}, dict
```

---

When there are multiple input variables, users can specify the name of the variable when initializing the instance. The name usually matches the variable name on the left hand side.

To get the derivative with respect to  $x$ , one simply uses "x" as the key to get the value from the dict `f.der`. Calling `f.der` gives the gradient with respect to all input variables inside a dictionary.

### 3.2.3 m inputs to n outputs

We recognize that users may sometimes want to compute Jacobians. Consider the following case,  $x = [x_1 \ x_2]^T, y = [y_1 \ y_2]^T, f(x, y) = [2x + \exp(y) \ 3x + 2 \sin(y)]^T$ .

We provide another class `Forward` to do the calculation.

---

```
variables = {'x': 3, 'y': 5} # specifies input values
functions = ['2*x + y', '3*x + 2*y'] # specifies functions
f = ad.Forward(variables, functions)
print(f)

# Output:
# Function: 2*x + y, Value: 11, Derivative: {'x': array(2.), 'y': array(1.)}
# Function: 3*x + 2*y, Value: 19, Derivative: {'x': array(3.), 'y': array(2.)}
```

---

In this case, all computed results will be saved in `f.results`, which is a list of values and derivatives of the two functions. The user can obtain the value and the derivatives of the  $i^{th}$  function by calling

---

```
val1 = f.results[0].val # 11
```

---

```
der1 = f.results[0].der # {'x': array(2.), 'y': array(1.)}  
val2 = f.results[1].val # 19  
der2 = f.results[1].der # {'x': array(3.), 'y': array(2.)}
```

---

## 4 Software Organization

### 4.1 Directory Structure

Our directory structure looks like the following:

```
cs107-FinalProject  
├── License  
├── README  
├── requirements.txt  
├── setup.py  
├── src/  
│   ├── __init__.py  
│   ├── utils.py  
│   ├── salad.py  
│   └── tests/  
│       ├── __init__.py  
│       ├── run_tests.sh  
│       └── test_forward.py  
└── docs/  
    ├── milestone1.pdf  
    ├── milestone2_progress.pdf  
    └── milestone2.pdf
```

### 4.2 Modules

Three external modules that we will rely on in our implementation will be NumPy, pytest, and coverage.

- NumPy offers support for large, multi-dimensional arrays as well a wide range of optimized mathematical functions and operators for these arrays. Numerical calculation of function values in Salad will rely extensively on this module.
- pytest is a feature rich testing framework in python that makes it easy to write unit tests that scale for full libraries.
- coverage is a commonly-used library used to compute code coverage during test execution.

### 4.3 Test Suite

The test suite will live inside the `tests` directory.

To run the test suite, find `run_tests.sh` in the `tests` folder. Run the following command to run test and generate coverage report.

---

```
sh run_tests.sh
```

---

Test coverage is currently at 96%.

## 4.4 Packaging and Distribution

We will eventually distribute the package through the Python Package Index (PyPi). Users will be able to install our package just by doing:

---

```
pip install salad
```

---

## 5 Implementation

### 5.1 Data Structure

The core data structures used in our implementation are `list`, `numpy.array`, and `dict`. See the “ $m$  inputs to  $n$  outputs” example from above as an illustration. The input variables are in a `dict`. The input functions are in a `dict`. The Jacobian with respect to each vector-form input is in `numpy.array`. To easily match the Jacobian to the corresponding input vector, we use a `dict`, where the key is the name of the corresponding input variable, and the value is a `numpy.array` of the Jacobian.

### 5.2 Classes and Methods

The baseline implementation currently includes two classes: 1) `Variable`, 2) `Forward`

#### 5.2.1 Variable

In `Variable`, we re-define dunder methods and other arithmetic operations. See below for the structure of the class and method definition.

---

```
class Variable(object):
    counter = 0

    def __init__(
        self, val, der=None, label=None, ad_mode="forward", increment_counter=True
    ):

        self.val = np.asarray(val)
        self.ad_mode = ad_mode

        if label is not None:
            self.label = label
        else:
            self.label = "v" + Variable.get_counter()

        if der is not None:
            self.der = der
        else:
            self.der = {self.label: np.ones(self.val.shape)}

        if increment_counter:
            Variable.increment()

    def __add__(self, other):
        pass
    def __str__(self):
        pass
    def __radd__(self, other):
```

```

    pass
def __sub__(self, other):
    pass
def __rsub__(self, other):
    pass
def __mul__(self, other):
    pass
def __rmul__(self, other):
    pass
def __truediv__(self, other):
    pass
def __rtruediv__(self, other):
    pass
def __neg__(self, other):
    pass
def __pow__(self, other):
    pass
def __rpow__(self, other):
    pass
def __eq__(self, other):
    pass
def __ne__(self, other):
    pass
def __lt__(self, other):
    pass
def __le__(self, other):
    pass
def __gt__(self, other):
    pass
def __ge__(self, other):
    pass

```

---

On a high level, each instance of this class corresponds to a node in the computation graph in the background section. We keep track of the primal trace (`self.val`) and the tangent trace (`self.der`) as attributes for all instances.

We also implement elementary functions like `sin`, `sqrt` (still need to be implemented), `log`, and `exp`.

```

def sin(x):
    pass
def cos(x):
    pass
def tan(x):
    pass
def exp(x):
    pass
def log(x):
    pass
def ln(x):
    pass
def sqrt{x}:
    pass
def logistic(x):
    pass

```

---

Users can use these functions for simple numeric calculations:

```

>>> x = 3
>>> ans = ad.sin(x)
>>> print(f'x = {ans}, \ntype of x: {type(ans)}')

```

```
x = 0.1411200080598672,  
type of x: <class 'numpy.float64'>
```

---

Or, users can use these functions to compute function derivatives. In this case, the inputs should be of the type `Variable`.

---

```
>>> x = ad.Variable(3, label="x")  
>>> ans = ad.sin(x)  
>>> print(f'x = {ans.val}, x.der = {ans.der}')  
x = 0.1411200080598672, x.der = {'x': -0.9899924966004454}  
>>> print(type(ans))  
<class 'salad.Variable'>
```

---

More arithmetic operations and functions to be included if time allows and as we see fit.

### 5.2.2 Forward

We also include a `Forward` class for easy handling of multiple function input. The variables are input in a `dict` and the functions are input in a `list`. The `__init__` function within `Forward` simply loops through and evaluate each functions one by one. See the “m inputs to n outputs” example from above as an illustration.

## 6 Licensing

We will use the MIT License. Future developers are welcome to extend our library and distribute closed source versions. Because we are using NumPy, a copyright licensed library, we are unable to use a license like GNU GPLv3 (as it requires all libraries used by the library to be copyleft).

## 7 Future Features

For our future features, we want to implement an optimization toolkit for our users. These optimization methods are useful in numerous contexts such as statistics, machine learning, operations research, and economics. The optimization methods we plan to implement are:

1. gradient descent
2. stochastic gradient descent
3. Broyden–Fletcher–Goldfarb–Shanno (BFGS)
4. Newton’s method (if time allows)

These functions will allow users to maximize or minimize an arbitrary differentiable function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ .

We note that the Newton’s method requires the Hessian, which requires a further extension to our implementation. A simple (but painful) fix is to ask users to supply the analytical first derivative(s), on which we perform the rooting finding iterations. However, ideally we would like to include the computation of the Hessian in our implementation, potentially as a third variable within the `Variable` class.

If time allows, we could compare and demonstrate the pros and cons of these optimization methods. For instance, stochastic gradient descent is more powerful than gradient descent at reaching the global minimum when the function is non-convex.

The implementation includes four different classes `GradientDescent`, `StochasticGradientDescent`, `BFGS`, `NewtonsMethod`, which inherit a baseline `Optimizer` class. All class definitions will be included in the `optimization.py` script in the `src` folder.