

Efficient differentiable programming

November 27, 2019

abstract

- an AD system for higher order functional array-processing language.
- supports both source-to-source forward mode AD and global optimizations
- gradient computation with forward AD can be as efficient as reverse mode
- Jacobian matrices needed for practical numerical schemes can be efficiently computed

The key contribution is in making traditional forward mode AD efficient through a number of compile time optimizing transformations, e.g. fusion

AD in general

Two main mode for AD:

- forward mode computes the derivative of the original computation while making a forward pass (across the computational graph)
- Reverse mode does a forward pass and then a backwards one to compute the adjoint derivatives.

For reverse mode we associate to every variable

$$v_i$$

its adjoint

$$\frac{\partial y}{\partial v_i}$$

We want a constant time computation of differentiation as an important assumption in algorithms is often that the gradient of a function has the same asymptotic computational cost as the original function.

Most algorithms in optimization have their computational complexity defined in the number of calls to an *oracle*. In **first order methods** the oracle is usually either gradient, a subgradient, or maybe some funky proximal operator.

supervised learning setup

AI is pretty simple nowadays, it's all machine learning, which means given a dataset of N samples of the form

$$\{(x_1, y_1), \dots, (x_N, y_N)\}$$

where the x 's are the *feature vectors* and the y 's the result of the prediction, we wish to find a good predictor, i.e. a function

$$g : X \rightarrow Y$$

To define the goodness of a predictor we define a *loss function*

$$f : Y \times Y \rightarrow \mathbb{R}$$

Define the empirical risk minimization has the average of all the losses over the sample. Then the best predictor is the one that minimizes the empirical risk.

- In general if

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

we define the Jacobian the following way:

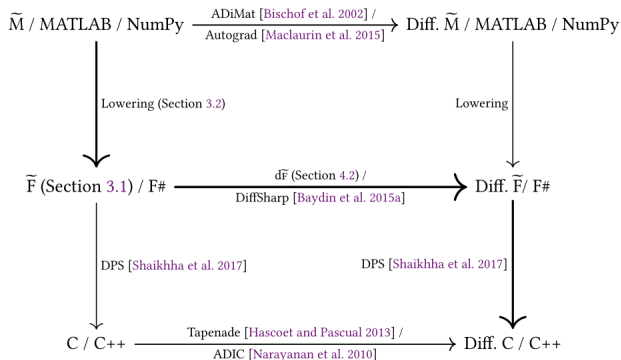
$$J_f = \frac{\partial f}{\partial x} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \dots & \frac{\partial f_m}{\partial x_n} \end{bmatrix} \quad (1)$$

Forward mode computes the Jacobian through columns whereas reverse mode computes rows.

Forward mode is efficient when there are many more rows than columns and vice versa for reverse mode. Reverse mode is thus good for very large dimensions in the input. Choosing the optimal mix between the two is NP complete.

- When computing the matrix there are many loop invariant computations, this calls for *loop-invariant code motion*
- Even if the asymptotic computational complexity of a gradient call is the same as one for the original function the hidden constant can be quite big in practice. Many intermediate vectors can be removed through **deforestation** aka **loop fusion**.

Pipeline



d f-smooth and M-smooth

- The highest level /api/ mimicks the conventional interfaces to high level linear algebra libraries and languages such as MatLab and Numpy
- The program is then lowered into its representation in f smooth which computes the actual derivatives. After this computation many optimizations can be applied before...
- Compiling once again to C code, where memory is managed in a stack like manner by using the destination passing style technique, which is just a fancy way to say that the caller is responsible for the callees memory and passes to it the memory location where it can store the results.

differentiation in df-smooth

- source to source transformation based on the `deriv` construct
 - everything done at compile time, macro-ish
 - we can call `deriv` multiple times to get the derivative with respect to a list of free variables.
- 1 first `deriv` constructs a lambda with the free variables as input parameters
 - 2 the produced lambda is passed to the source-to-source differentiation operator \mathcal{D}

The derivative of the program `cos a` is represented as `snd((\mathcal{D} [[fun a -> cos a]])(a, 1))` which expands to

$$\text{snd}((\mathcal{D}[\text{fun } a \rightarrow \cos a])(a, 1))$$

Example 2 (Continued). In the previous example, based on the automatic differentiation rules, the differentiated program would be as follows:

```
 $\vec{g} = \text{fun } \vec{x} \rightarrow -\text{snd } (\vec{x}) * \sin(\text{fst } (\vec{x}))$ 
```

Based on the definition of the `diff` construct, we have to use the AD version of the function (i.e., `g`) and assign 1 to the derivative part of the input. So the value of `cos'` for the input `a` is computed as follows:

$$\begin{aligned} \text{snd } ((\text{diff } g) \ a) &\rightsquigarrow \text{snd } (\mathcal{D}[\vec{g}](a, 1)) \rightsquigarrow \text{snd } (\vec{g} \ (a, 1)) \rightsquigarrow \\ &-\text{snd } ((a, 1)) * \sin(\text{fst } ((a, 1))) \rightsquigarrow -1 * \sin(a) \rightsquigarrow -\sin(a) \end{aligned}$$

△

$e ::=$	$e \bar{e} \mid \text{fun } \bar{x} \rightarrow e \mid x$	- Application, Abstraction, and Variable Access
	$ n \mid i \mid N$	- Scalar, Index, and Cardinality Value
	$ c$	- Constants (see below)
	$ \text{let } x = e \text{ in } e$	- (Non-Recursive) Let Binding
	$ \text{if } e \text{ then } e \text{ else } e$	- Conditional
$T ::=$	M	- (Non-Functional) Expression Type
	$ \bar{T} \Rightarrow M$	- Function Types (No Currying)
$M ::=$	Num	- Numeric Type
	$ \text{Array}\langle M \rangle$	- Vector, Matrix, ... Type
	$ M \times M$	- Pair Type
	$ \text{Bool}$	- Boolean Type
$\text{Num} ::=$	$\text{Double} \mid \text{Index} \mid \text{Card}$	- Scalar, Index, and Cardinality Type

Typing Rules:

$$\begin{array}{c}
\text{(T-App)} \frac{e_0 : \bar{T} \Rightarrow M \quad \bar{e} : \bar{T}}{e_0 \bar{e} : M} \quad \text{(T-Abs)} \frac{\Gamma \cup \bar{x} : \bar{T} \vdash e : M}{\Gamma \vdash \lambda \bar{x}. e : \bar{T} \Rightarrow M} \quad \text{(T-Var)} \frac{x : T \in \Gamma}{\Gamma \vdash x : T} \\
\text{(T-Let)} \frac{\Gamma \vdash e_1 : T_1 \quad \Gamma, x : T_1 \vdash e_2 : T_2}{\Gamma \vdash \text{let } x = e_1 \text{ in } e_2 : T_2} \quad \text{(T-If)} \frac{e_1 : \text{Bool} \quad e_2 : M \quad e_3 : M}{\text{if } e_1 \text{ then } e_2 \text{ else } e_3 : M}
\end{array}$$

Vector Function Constants:

$\text{build} : \text{Card} \Rightarrow (\text{Index} \Rightarrow M) \Rightarrow \text{Array}\langle M \rangle$ $\text{get} : \text{Array}\langle M \rangle \Rightarrow \text{Index} \Rightarrow M$
 $\text{ifold} : (M \Rightarrow \text{Index} \Rightarrow M) \Rightarrow M \Rightarrow \text{Card} \Rightarrow M$ $\text{length} : \text{Array}\langle M \rangle \Rightarrow \text{Card}$

Pair Function Constants:

$\text{pair} : M_1 \Rightarrow M_2 \Rightarrow M_1 \times M_2$ $\text{fst} : M_1 \times M_2 \Rightarrow M_1$ $\text{snd} : M_1 \times M_2 \Rightarrow M_2$

Syntactic Sugar:

$e_0[e_1]$	=	$\text{get } e_0 \ e_1$	Matrix	=	$\text{Array}\langle \text{Array}\langle \text{Double} \rangle \rangle$
(e_0, e_1)	=	$\text{pair } e_0 \ e_1$	DoubleD	=	$\text{Double} \times \text{Double}$
$e_1 \text{ bop } e_2$	=	$\text{bop } e_1 \ e_2$	VectorD	=	$\text{Array}\langle \text{Double} \times \text{Double} \rangle$
Vector	=	$\text{Array}\langle \text{Double} \rangle$	MatrixD	=	$\text{Array}\langle \text{Array}\langle \text{Double} \times \text{Double} \rangle \rangle$

Fig. 3. The syntax, type system, and function constants of the core \tilde{F} .

built in functions for array programming : build ifold and length

loop fusion example

In this example we show that taking the tranpose twice of a Matrix is the “same” as just copying it.

Efficient Differentiable Programming in a Functional Array-Processing Language

97:15

```
let MT =  
  build (length M[0]) (fun i ->  
    build (length M) (fun j ->  
      M[j][i] ) ) in  
  build (length MT[0]) (fun i ->  
    build (length MT) (fun j ->  
      MT[j][i] ) )
```

Now, by applying the loop fusion rules (cf. Figure 8c) and performing further partial evaluation, the following expression is derived:

```
build (length M) (fun i ->  
  build (length M[0]) (fun j ->  
    M[i][j] ) )
```

This is the same expression as M.

- standard ml compilers/libraries : pytorch, tensorflow
- stalingrad
- Pilatus
- category approach
- other numerical DSL: Halide

The nitty gritty details of AD can be horrendous in terms of compiler optimization etc. The authors show how complicated it can be to represent differentiation at a high level and still keep the implementation efficient.

However the authors should have made a choice between high level and compiler optimizations. Elliot's approach is "wayy" more general, the high level api of df smooth is mostly there to enforce invariants to enable optimization.

The question of how to effectively model mathematical properties which are simple for cegep students but very hard to encode precisely is more important.