

# Taking Derivatives of Functional Programs

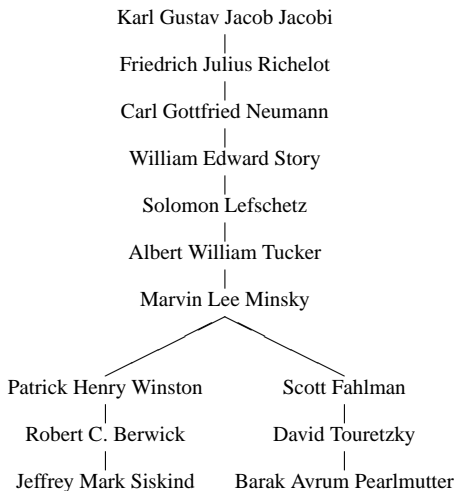
## AD in a Functional Framework

Jeffrey Mark Siskind  
qobi@purdue.edu

School of Electrical and Computer Engineering  
Purdue University

New England Programming Languages and Systems  
27 October 2005

Joint work with Barak A. Pearlmutter.



# Outline

- 1 Introduction
- 2 Tutorial on AD
  - Forward Mode
  - Reverse Mode
- 3 Examples

# Outline

- 1 Introduction
- 2 Tutorial on AD
  - Forward Mode
  - Reverse Mode
- 3 Examples

*It is, of course, not excluded that the range of arguments or range of values of a function should consist wholly or partly of functions. The derivative, as this notion appears in the elementary differential calculus, is a familiar mathematical example of a function for which both ranges consist of functions.*

(p. 1 ¶4)

Church, A. (1941). *The Calculi of Lambda Conversion*. Princeton University Press, Princeton, NJ.

Gottfried Leibniz  
|  
Jacob Bernoulli  
|  
Johann Bernoulli  
|  
Leonhard Euler  
|  
Joseph Louis Lagrange  
|  
Simeon Poisson  
|  
Michel Chasles  
|  
Hubert Anson Newton  
|  
Eliakim Hastings Moore  
|  
Oswald Veblen  
|  
Alonzo Church

Leibnitz (1664) + Church (1941) = Siskind & Pearlmutter (2005)

Leibnitz, G. W. (1664). A new method for maxima and minima as well as tangents, which is impeded neither by fractional nor irrational quantities, and a remarkable type of calculus for this. *Acta Eruditorum*.

# Differential Calculus for Dummies

(in 6 slides)

# Derivatives

$$\frac{dax^2}{dx} \rightsquigarrow 2ax$$

# Derivatives

$$\frac{dax^2}{dx} \rightsquigarrow 2ax$$

$$\frac{d}{dx} : \underbrace{f}_{\mathbb{R} \rightarrow \mathbb{R}} \mapsto \underbrace{f'}_{\mathbb{R} \rightarrow \mathbb{R}}$$

# Derivatives

$$\frac{dax^2}{dx} \rightsquigarrow 2ax$$

$$\frac{d}{dx} : \underbrace{f}_{\mathbb{R} \rightarrow \mathbb{R}} \mapsto \underbrace{f'}_{\mathbb{R} \rightarrow \mathbb{R}}$$

$$\frac{d}{dx} : (\mathbb{R} \rightarrow \mathbb{R}) \rightarrow (\mathbb{R} \rightarrow \mathbb{R})$$

# Derivatives

$$\frac{dax^2}{dx} \rightsquigarrow 2ax$$

$$\frac{d}{dx} : \underbrace{f}_{\mathbb{R} \rightarrow \mathbb{R}} \mapsto \underbrace{f'}_{\mathbb{R} \rightarrow \mathbb{R}}$$

$$\frac{d}{dx} : (\mathbb{R} \rightarrow \mathbb{R}) \rightarrow (\mathbb{R} \rightarrow \mathbb{R})$$

$$\mathcal{D} : (\mathbb{R} \rightarrow \mathbb{R}) \rightarrow (\mathbb{R} \rightarrow \mathbb{R})$$

# Derivatives

$$\frac{dax^2}{dx} \rightsquigarrow 2ax$$

$$\frac{d}{dx} : \underbrace{f}_{\mathbb{R} \rightarrow \mathbb{R}} \mapsto \underbrace{f'}_{\mathbb{R} \rightarrow \mathbb{R}}$$

$$\frac{d}{dx} : (\mathbb{R} \rightarrow \mathbb{R}) \rightarrow (\mathbb{R} \rightarrow \mathbb{R})$$

$$\mathcal{D} : (\mathbb{R} \rightarrow \mathbb{R}) \rightarrow (\mathbb{R} \rightarrow \mathbb{R})$$

$$\mathcal{D} \lambda x ax^2$$

# Partial Derivatives

$$\frac{\partial ax^2y^3}{\partial x}$$

$$\frac{\partial ax^2y^3}{\partial y}$$

# Partial Derivatives

$$\frac{\partial ax^2y^3}{\partial x}$$

$$\mathcal{D} \lambda x ax^2y^3$$

$$\frac{\partial ax^2y^3}{\partial y}$$

$$\mathcal{D} \lambda y ax^2y^3$$

# Partial Derivatives

$$\frac{\partial ax^2y^3}{\partial x}$$

$$\mathcal{D} \lambda x \ ax^2y^3$$

$$\mathcal{D}_1 \ \lambda(x, y) \ ax^2y^3$$

$$\frac{\partial ax^2y^3}{\partial y}$$

$$\mathcal{D} \ \lambda y \ ax^2y^3$$

$$\mathcal{D}_2 \ \lambda(x, y) \ ax^2y^3$$

# Partial Derivatives

$$\frac{\partial ax^2y^3}{\partial x}$$

$$\frac{\partial ax^2y^3}{\partial y}$$

$$\mathcal{D} \lambda x ax^2y^3$$

$$\mathcal{D} \lambda y ax^2y^3$$

$$\mathcal{D}_1 \lambda(x, y) ax^2y^3$$

$$\mathcal{D}_2 \lambda(x, y) ax^2y^3$$

$$\frac{\partial}{\partial x} : \underbrace{f}_{\mathbb{R}^n \rightarrow \mathbb{R}} \mapsto \underbrace{f'}_{\mathbb{R}^n \rightarrow \mathbb{R}}$$

# Partial Derivatives

$$\frac{\partial ax^2y^3}{\partial x}$$

$$\frac{\partial ax^2y^3}{\partial y}$$

$$\mathcal{D} \lambda x ax^2y^3$$

$$\mathcal{D} \lambda y ax^2y^3$$

$$\mathcal{D}_1 \lambda(x, y) ax^2y^3$$

$$\mathcal{D}_2 \lambda(x, y) ax^2y^3$$

$$\frac{\partial}{\partial x} : \underbrace{f}_{\mathbb{R}^n \rightarrow \mathbb{R}} \mapsto \underbrace{f'}_{\mathbb{R}^n \rightarrow \mathbb{R}}$$

$$\frac{\partial}{\partial x} : (\mathbb{R}^n \rightarrow \mathbb{R}) \rightarrow (\mathbb{R}^n \rightarrow \mathbb{R})$$

# Partial Derivatives

$$\frac{\partial ax^2y^3}{\partial x}$$

$$\frac{\partial ax^2y^3}{\partial y}$$

$$\mathcal{D} \lambda x ax^2y^3$$

$$\mathcal{D} \lambda y ax^2y^3$$

$$\mathcal{D}_1 \lambda(x, y) ax^2y^3$$

$$\mathcal{D}_2 \lambda(x, y) ax^2y^3$$

$$\frac{\partial}{\partial x} : \underbrace{f}_{\mathbb{R}^n \rightarrow \mathbb{R}} \mapsto \underbrace{f'}_{\mathbb{R}^n \rightarrow \mathbb{R}}$$

$$\frac{\partial}{\partial x} : (\mathbb{R}^n \rightarrow \mathbb{R}) \rightarrow (\mathbb{R}^n \rightarrow \mathbb{R})$$

$$\mathcal{D}_i : (\mathbb{R}^n \rightarrow \mathbb{R}) \rightarrow (\mathbb{R}^n \rightarrow \mathbb{R})$$

# Gradients

$$\nabla f \mathbf{x} = (\mathcal{D}_1 f \mathbf{x}), \dots, (\mathcal{D}_n f \mathbf{x})$$

$$\nabla : (\mathbb{R}^n \rightarrow \mathbb{R}) \rightarrow (\mathbb{R}^n \rightarrow \mathbb{R}^n)$$

# Jacobians

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

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

$$(\mathcal{J} f \mathbf{x})[i,j] = (\nabla (\mathbf{f}[i]))[j]$$

$$\mathcal{J} : (\mathbb{R}^n \rightarrow \mathbb{R}^m) \rightarrow (\mathbb{R}^n \rightarrow \mathbb{R}^{m \times n})$$

# Operators

$\mathcal{D}$ ,  $\nabla$ , and  $\mathcal{J}$  are traditionally called *operators*.

A more modern term is *higher-order functions*.

Higher-order functions are common in mathematics, physics, and engineering:

*summations, comprehensions, quantifications, optimizations, integrals, convolutions, filters, edge detectors, Fourier transforms, differential equations, Hamiltonians, . . .*

# The Chain Rule

$$(f \circ g) x = g (f x)$$

# The Chain Rule

$$(f \circ g) x = g (f x)$$

$$\frac{dg}{dx} = \frac{dg}{df} \frac{df}{dx}$$

# The Chain Rule

$$(f \circ g) x = g (f x)$$

$$\frac{dg}{dx} = \frac{dg}{df} \frac{df}{dx}$$

$$\mathcal{D} (f \circ g) x = (\mathcal{D} g (f x)) \times (\mathcal{D} f x)$$

# The Chain Rule

$$(f \circ g) x = g (f x)$$

$$\frac{dg}{dx} = \frac{dg}{df} \frac{df}{dx}$$

$$\mathcal{D} (f \circ g) x = (\mathcal{D} g (f x)) \times (\mathcal{D} f x)$$

$$\mathcal{J} (f \circ g) \mathbf{x} = (\mathcal{J} g (f \mathbf{x})) \times (\mathcal{J} f \mathbf{x})$$

# Outline

- 1 Introduction
- 2 Tutorial on AD**
  - Forward Mode
  - Reverse Mode
- 3 Examples

# Outline

1 Introduction

2 Tutorial on AD

- Forward Mode
- Reverse Mode

3 Examples

# Straight-Line Code and Jacobians

$$\mathbf{x}_1 = f_1 \mathbf{x}_0$$

$$\vdots$$

$$\mathbf{x}_n = f_n \mathbf{x}_{n-1}$$

$$f = f_1 \circ \cdots \circ f_n$$

$$\mathcal{J} f \mathbf{x}_0 = (\mathcal{J} f_n \mathbf{x}_{n-1}) \times \cdots \times (\mathcal{J} f_1 \mathbf{x}_0)$$

$$(\mathcal{J} f \mathbf{x}_0)^\top = (\mathcal{J} f_1 \mathbf{x}_0)^\top \times \cdots \times (\mathcal{J} f_n \mathbf{x}_{n-1})^\top$$

# One Way to Compute the Jacobian

$$\overline{\mathbf{X}}_1' = (\mathcal{J} f_1 \mathbf{x}_0)$$

$$\overline{\mathbf{X}}_2' = (\mathcal{J} f_2 \mathbf{x}_1) \times \overline{\mathbf{X}}_1'$$

$$\vdots$$

$$\overline{\mathbf{X}}_n' = (\mathcal{J} f_n \mathbf{x}_{n-1}) \times \overline{\mathbf{X}}_{n-1}'$$

$$\overline{\mathbf{X}}_n' = \mathcal{J} f \mathbf{x}_0$$

# Forward-Mode AD

$$\overline{\mathbf{x}}_1' = (\mathcal{J} f_1 \mathbf{x}_0) \times \overline{\mathbf{x}}_0'$$

$$\vdots$$

$$\overline{\mathbf{x}}_n' = (\mathcal{J} f_n \mathbf{x}_{n-1}) \times \overline{\mathbf{x}}_{n-1}'$$

$$\overline{\mathbf{x}}_n' = (\mathcal{J} f \mathbf{x}_0) \times \overline{\mathbf{x}}_0'$$

Wengert, R. E. (1964). A simple automatic derivative evaluation program. *Communications of the ACM*, **7**(8):463–4.

# Interleaving Forward Mode

$$\mathbf{x}_1 = f_1 \mathbf{x}_0$$

$$\vdots$$

$$\mathbf{x}_n = f_n \mathbf{x}_{n-1}$$

$$\overline{\mathbf{x}}_1 = (\mathcal{J} f_1 \mathbf{x}_0) \times \overline{\mathbf{x}}_0$$

$$\vdots$$

$$\overline{\mathbf{x}}_n = (\mathcal{J} f_n \mathbf{x}_{n-1}) \times \overline{\mathbf{x}}_{n-1}$$

$$\mathbf{x}_1 = f_1 \mathbf{x}_0$$

$$\overline{\mathbf{x}}_1 = (\mathcal{J} f_1 \mathbf{x}_0) \times \overline{\mathbf{x}}_0$$

$$\vdots$$

$$\mathbf{x}_n = f_n \mathbf{x}_{n-1}$$

$$\overline{\mathbf{x}}_n = (\mathcal{J} f_n \mathbf{x}_{n-1}) \times \overline{\mathbf{x}}_{n-1}$$

# Forward Mode as a Transformation

$$\left. \begin{array}{l} \mathbf{x}_1 = f_1 \mathbf{x}_0 \\ \vdots \\ \mathbf{x}_n = f_n \mathbf{x}_{n-1} \end{array} \right\} \rightsquigarrow \left\{ \begin{array}{l} \overrightarrow{\mathbf{x}}_1 = \overrightarrow{f}_1 \overrightarrow{\mathbf{x}}_0 \\ \vdots \\ \overrightarrow{\mathbf{x}}_n = \overrightarrow{f}_n \overrightarrow{\mathbf{x}}_{n-1} \end{array} \right.$$

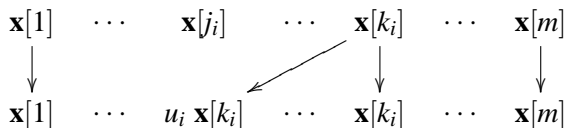
$$\begin{aligned} \overrightarrow{\mathbf{x}} &= (\mathbf{x}, \overline{\mathbf{x}}) \\ \overrightarrow{f}(\mathbf{x}, \overline{\mathbf{x}}) &= ((f \mathbf{x}), ((\mathcal{J} f \mathbf{x}) \times \overline{\mathbf{x}})) \end{aligned}$$

# A Unary Sparse Function

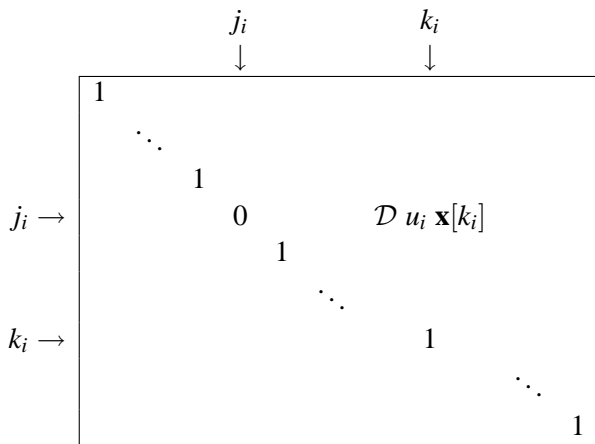
$$(f_i \mathbf{x})[j_i] = u_i \mathbf{x}[k_i]$$

$$(f_i \mathbf{x})[j'] = \mathbf{x}[j']$$

$$j' \neq j_i$$



# The Jacobian of a Unary Sparse Function



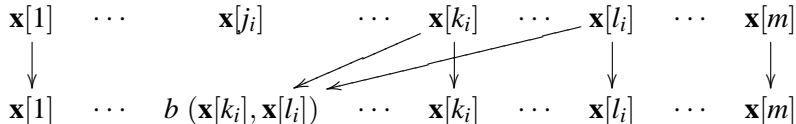


# A Binary Sparse Function

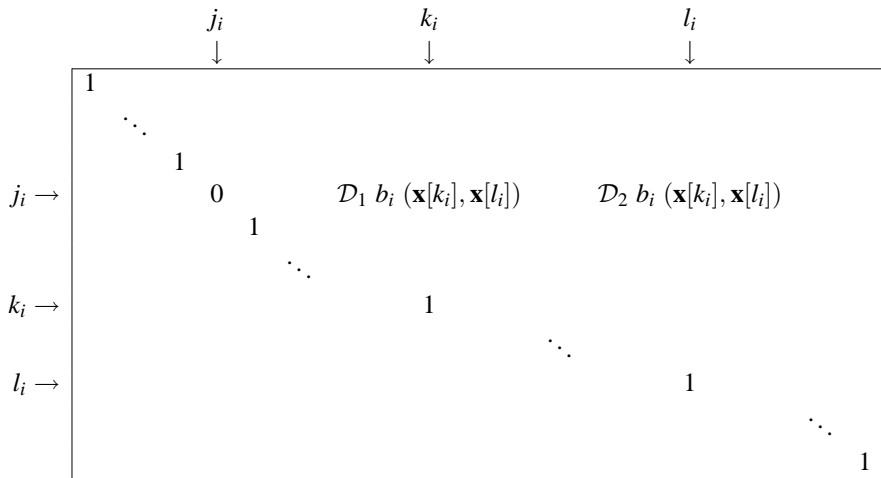
$$(f_i \mathbf{x})[j_i] = b_i(\mathbf{x}[k_i], \mathbf{x}[l_i])$$

$$(f_i \mathbf{x})[j'] = \mathbf{x}[j']$$

$$j' \neq j_i$$



# The Jacobian of a Binary Sparse Function





# Forward Mode as a Sparse Transformation

$$\begin{aligned}
 x_{j_i} &:= u_i x_{k_i} &\rightsquigarrow& \overrightarrow{x_{j_i}} := \overrightarrow{u_i} \overrightarrow{x_{k_i}} \\
 x_{j_i} &:= b_i (x_{k_i}, x_{l_i}) &\rightsquigarrow& \overrightarrow{x_{j_i}} := \overrightarrow{b_i} (\overrightarrow{x_{k_i}}, \overrightarrow{x_{l_i}})
 \end{aligned}$$

$$\begin{aligned}
 \overrightarrow{x} &= (x, \overline{x'}) \\
 \overrightarrow{u} (x, \overline{x'}) &= ((u x), ((\mathcal{D} u x) \times \overline{x'})) \\
 \overrightarrow{b} ((x_1, \overline{x'_1}), (x_2, \overline{x'_2})) &= ((b (x_1, x_2)), \\
 &\quad (((\mathcal{D}_1 b (x_1, x_2)) \times \overline{x'_1}) + ((\mathcal{D}_2 b (x_1, x_2)) \times \overline{x'_2})))
 \end{aligned}$$

# Outline

- 1 Introduction
- 2 **Tutorial on AD**
  - Forward Mode
  - **Reverse Mode**
- 3 Examples

# Straight-Line Code and Jacobians

$$\mathbf{x}_1 = f_1 \mathbf{x}_0$$

$$\vdots$$

$$\mathbf{x}_n = f_n \mathbf{x}_{n-1}$$

$$f = f_1 \circ \cdots \circ f_n$$

$$\mathcal{J} f \mathbf{x}_0 = (\mathcal{J} f_n \mathbf{x}_{n-1}) \times \cdots \times (\mathcal{J} f_1 \mathbf{x}_0)$$

$$(\mathcal{J} f \mathbf{x}_0)^\top = (\mathcal{J} f_1 \mathbf{x}_0)^\top \times \cdots \times (\mathcal{J} f_n \mathbf{x}_{n-1})^\top$$

# Another Way to Compute the Jacobian

$$\overline{\mathbf{x}}_{n-1} = (\mathcal{J} f_n \mathbf{x}_{n-1})^\top$$

$$\overline{\mathbf{x}}_{n-2} = (\mathcal{J} f_{n-1} \mathbf{x}_{n-2})^\top \times \overline{\mathbf{x}}_{n-1}$$

$$\vdots$$

$$\overline{\mathbf{x}}_0 = (\mathcal{J} f_1 \mathbf{x}_0)^\top \times \overline{\mathbf{x}}_1$$

$$\overline{\mathbf{x}}_0 = (\mathcal{J} f \mathbf{x}_0)^\top$$

# Reverse-Mode AD

$$\overline{\mathbf{x}}_{n-1} = (\mathcal{J} f_n \mathbf{x}_{n-1})^\top \times \overline{\mathbf{x}}_n$$

$$\vdots$$

$$\overline{\mathbf{x}}_0 = (\mathcal{J} f_1 \mathbf{x}_0)^\top \times \overline{\mathbf{x}}_1$$

$$\overline{\mathbf{x}}_0 = (\mathcal{J} f \mathbf{x}_0)^\top \times \overline{\mathbf{x}}_n$$

Speelpenning, B. (1980). *Compiling Fast Partial Derivatives of Functions Given by Algorithms*. PhD thesis, Department of Computer Science, University of Illinois at Urbana-Champaign.

# Reverse Mode Cannot be Interleaved

$$\mathbf{x}_1 = f_1 \mathbf{x}_0$$

$$\vdots$$

$$\mathbf{x}_n = f_n \mathbf{x}_{n-1}$$

$$\overline{\mathbf{x}}_{n-1} = (\mathcal{J} f_n \mathbf{x}_{n-1})^\top \times \overline{\mathbf{x}}_n$$

$$\vdots$$

$$\overline{\mathbf{x}}_0 = (\mathcal{J} f_1 \mathbf{x}_0)^\top \times \overline{\mathbf{x}}_1$$

# Reverse Mode via Backpropagators

$$\mathbf{x}_1 = f_1 \mathbf{x}_0$$

$$\overline{\mathbf{x}}_1 = \lambda \overline{\mathbf{x}} \overline{\mathbf{x}}_0 ((\mathcal{J} f_1 \mathbf{x}_0)^\top \times \overline{\mathbf{x}})$$

$$\vdots$$

$$\mathbf{x}_n = f_n \mathbf{x}_{n-1}$$

$$\overline{\mathbf{x}}_n = \lambda \overline{\mathbf{x}} \overline{\mathbf{x}}_{n-1} ((\mathcal{J} f_n \mathbf{x}_{n-1})^\top \times \overline{\mathbf{x}})$$

$$\overline{\mathbf{x}}_n \overline{\mathbf{x}}_n$$

# Reverse Mode as a Transformation

$$\left. \begin{array}{l} \mathbf{x}_1 = f_1 \mathbf{x}_0 \\ \vdots \\ \mathbf{x}_n = f_n \mathbf{x}_{n-1} \end{array} \right\} \rightsquigarrow \left\{ \begin{array}{l} \overleftarrow{\mathbf{x}}_1 = \overleftarrow{f}_1 \overleftarrow{\mathbf{x}}_0 \\ \vdots \\ \overleftarrow{\mathbf{x}}_n = \overleftarrow{f}_n \overleftarrow{\mathbf{x}}_{n-1} \end{array} \right.$$

$$\begin{aligned} \overleftarrow{\mathbf{x}} &= (\mathbf{x}, \bar{\mathbf{x}}) \\ \overleftarrow{f}(\mathbf{x}, \bar{\mathbf{x}}) &= ((f \mathbf{x}), (\lambda \bar{\mathbf{x}} \bar{\mathbf{x}} ((\mathcal{J} f \mathbf{x})^\top \times \bar{\mathbf{x}}))) \end{aligned}$$

# Reverse Mode via a Tape

$$\left. \begin{array}{l} \mathbf{x}_1 = f_1 \mathbf{x}_0 \\ \vdots \\ \mathbf{x}_n = f_n \mathbf{x}_{n-1} \end{array} \right\} \rightsquigarrow \left\{ \begin{array}{l} \overleftarrow{\mathbf{x}}_1 = \overleftarrow{f}_1 \overleftarrow{\mathbf{x}}_0 \\ \vdots \\ \overleftarrow{\mathbf{x}}_n = \overleftarrow{f}_n \overleftarrow{\mathbf{x}}_{n-1} \end{array} \right.$$

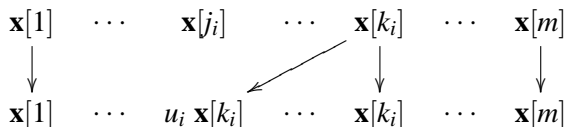
$$\begin{aligned} \overleftarrow{\mathbf{x}} &= \mathbf{x} \\ \overleftarrow{f} \mathbf{x} &= \mathbf{begin} \overline{\mathbf{x}} := \lambda \overline{\mathbf{x}} \overline{\mathbf{x}} ((\mathcal{J} f \mathbf{x})^\top \times \overline{\mathbf{x}}); \\ &\quad (f \mathbf{x}) \mathbf{end} \end{aligned}$$

# A Unary Sparse Function

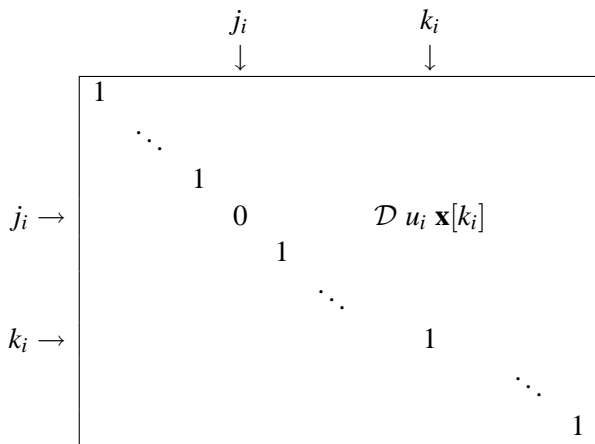
$$(f_i \mathbf{x})[j_i] = u_i \mathbf{x}[k_i]$$

$$(f_i \mathbf{x})[j'] = \mathbf{x}[j']$$

$$j' \neq j_i$$



# The Jacobian of a Unary Sparse Function





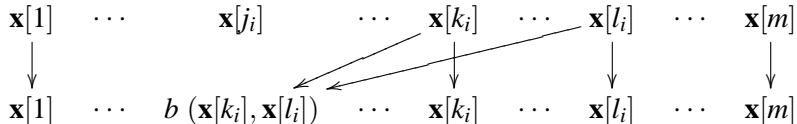


# A Binary Sparse Function

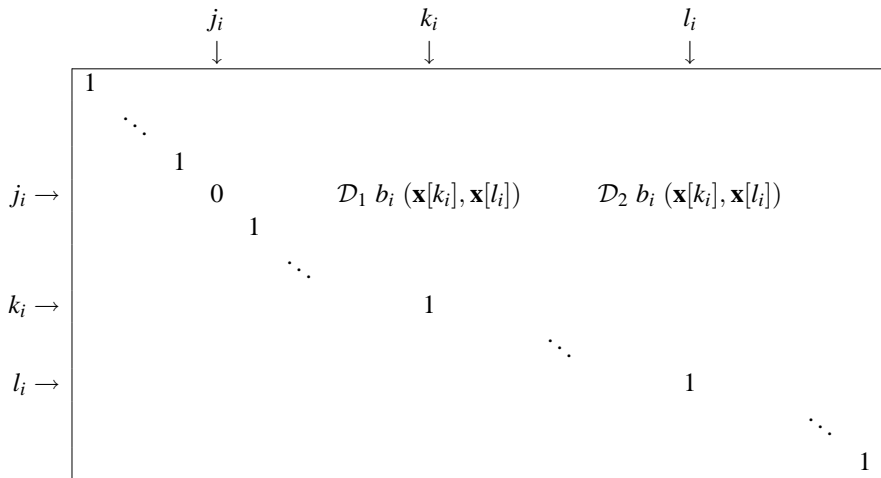
$$(f_i \mathbf{x})[j_i] = b_i(\mathbf{x}[k_i], \mathbf{x}[l_i])$$

$$(f_i \mathbf{x})[j'] = \mathbf{x}[j']$$

$$j' \neq j_i$$



# The Jacobian of a Binary Sparse Function







# Sparse Reverse Mode via a Tape

$$\begin{aligned}
 x_{j_i} := u_i x_{k_i} &\rightsquigarrow \bar{x} := \lambda[ ] \mathbf{begin} \quad \overline{x_{k_i}} += (\mathcal{D} u_i x_{k_i}) \times \overline{x_{j_i}}; \\
 &\quad \overline{x_{j_i}} := 0; \\
 &\quad \bar{x} [ ] \mathbf{end}; \\
 & \\
 &x_{j_i} := u_i x_{k_i}
 \end{aligned}$$

$$\begin{aligned}
 x_{j_i} := b_i(x_{k_i}, x_{l_i}) &\rightsquigarrow \bar{x} := \lambda[ ] \mathbf{begin} \quad \overline{x_{k_i}} += (\mathcal{D}_1 b_i(x_{k_i}, x_{l_i})) \times \overline{x_{j_i}}; \\
 &\quad \overline{x_{l_i}} += (\mathcal{D}_2 b_i(x_{k_i}, x_{l_i})) \times \overline{x_{j_i}}; \\
 &\quad \overline{x_{j_i}} := 0; \\
 &\quad \bar{x} [ ] \mathbf{end}; \\
 & \\
 &x_{j_i} := b_i(x_{k_i}, x_{l_i})
 \end{aligned}$$

# Traditional AD

Forward Mode:  $\mathbb{R}^n \rightarrow \mathbb{R}^m \rightsquigarrow (\mathbb{R}^n \times \mathbb{R}^n) \rightarrow (\mathbb{R}^m \times \mathbb{R}^m)$

Reverse Mode:  $\mathbb{R}^n \rightarrow \mathbb{R}^m \rightsquigarrow \mathbb{R}^n \rightarrow (\mathbb{R}^m \times (\mathbb{R}^m \rightarrow \mathbb{R}^n))$

## Functional AD

$$\begin{aligned} \overline{\text{null}} &= \text{null} \\ \overline{\mathbb{R}} &= \mathbb{R} \\ \overline{\tau_1 \times \tau_2} &= \overline{\tau_1} \times \overline{\tau_2} \\ \overline{\tau_1 \xrightarrow{\tau'_1, \dots, \tau'_n} \tau_2} &= \overline{\tau_1} \times \dots \times \overline{\tau_n} \end{aligned}$$

$$\begin{aligned} \overrightarrow{\text{null}} &= \text{null} \times \overline{\text{null}} \\ \overrightarrow{\mathbb{R}} &= \mathbb{R} \times \overline{\mathbb{R}} \\ \overrightarrow{\tau_1 \times \tau_2} &= \overrightarrow{\tau_1} \times \overrightarrow{\tau_2} \\ \overrightarrow{\tau_1 \xrightarrow{\tau'_1, \dots, \tau'_n} \tau_2} &= \overrightarrow{\tau_1} \xrightarrow{\tau'_1, \dots, \tau'_n} \overrightarrow{\tau_2} \end{aligned}$$

$$\begin{aligned} \overleftarrow{\text{null}} &= \text{null} \\ \overleftarrow{\mathbb{R}} &= \mathbb{R} \\ \overleftarrow{\tau_1 \times \tau_2} &= \overleftarrow{\tau_1} \times \overleftarrow{\tau_2} \\ \overleftarrow{\tau_1 \xrightarrow{\tau'_1, \dots, \tau'_n} \tau_2} &= \overleftarrow{\tau_1} \times \dots \times \overleftarrow{\tau_n} \end{aligned}$$

$$\begin{aligned} \overleftarrow{\text{null}} &= \text{null} \\ \overleftarrow{\mathbb{R}} &= \mathbb{R} \\ \overleftarrow{\tau_1 \times \tau_2} &= \overleftarrow{\tau_1} \times \overleftarrow{\tau_2} \\ \overleftarrow{\tau_1 \xrightarrow{\tau'_1, \dots, \tau'_n} \tau_2} &= \\ \overleftarrow{\tau_1} \xrightarrow{\overleftarrow{\tau'_1}, \dots, \overleftarrow{\tau'_n}} & (\overleftarrow{\tau_2} \times (\overleftarrow{\tau_2} \rightarrow (\overleftarrow{\tau'_1} \times \dots \times \overleftarrow{\tau'_n}) \times \overleftarrow{\tau_1})) \end{aligned}$$

Forward Mode:  $\overrightarrow{\mathcal{J}} : \tau \rightarrow \overrightarrow{\tau}$

Reverse Mode:  $\overleftarrow{\mathcal{J}} : \tau \rightarrow \overleftarrow{\tau}$

# Outline

- 1 Introduction
- 2 Tutorial on AD
  - Forward Mode
  - Reverse Mode
- 3 Examples

# Derivatives

$$\mathcal{D} f x \triangleq \text{TANGENT } ((\overrightarrow{\mathcal{J}} f) (x \blacktriangleright 1))$$

$$\mathcal{D} f x \triangleq \text{CDR } ((\text{CDR } ((\overleftarrow{\mathcal{J}} f) (\overleftarrow{\mathcal{J}} x)))) 1)$$

# Roots using Newton-Raphson

$$\text{ROOT } (f, x_0, \epsilon) \triangleq \mathbf{let } x' \triangleq x_0 - \frac{f x_0}{\mathcal{D} f x_0}$$

$$\mathbf{in if } |x_0 - x'| \leq \epsilon \mathbf{ then } x_0 \mathbf{ else } \text{ROOT } (f, x', \epsilon)$$

# Univariate Minimizer

## Line Search

$$\text{LINESEARCH}(f, x_0, \epsilon) \triangleq \text{ROOT}((\mathcal{D}f), x_0, \epsilon)$$

## Gradients

$$\nabla f x \triangleq \mathbf{let} \ n \triangleq \mathbf{LENGTH} \ x$$

$$\mathbf{in} \ \mathbf{MAP} \ ((\lambda i \ \mathbf{TANGENT} \ ((\vec{\mathcal{J}} f) (x \blacktriangleright e_{i,n}))), (\iota n))$$

$$\nabla f x \triangleq \mathbf{CDR} \ ((\mathbf{CDR} \ ((\overleftarrow{\mathcal{J}} f) (\overleftarrow{\mathcal{J}} x))) \ 1)$$

# Multivariate Minimizer

## Gradient Descent

```

GRADIENDESCENT ( $f, x_0, \epsilon$ )  $\triangleq$ 
  let  $g \triangleq \nabla f x_0$ 
  in if  $\|g\| \leq \epsilon$ 
    then  $x_0$ 
    else GRADIENDESCENT
      ( $f, (x_0 + ((\text{LINESEARCH } ((\lambda k f (x_0 + (k \times g))), \epsilon)) \times g)), \epsilon$ )
  
```

# Saddle Points

## Continuous Two-Person Zero Sum Games

$$\mathbf{x} : \mathbb{R}^m$$

$$\mathbf{y} : \mathbb{R}^n$$

$$\text{PAYOFF} : \mathbb{R}^m \times \mathbb{R}^n \rightarrow \mathbb{R}$$

$$\min_{\mathbf{x}} \max_{\mathbf{y}} \text{PAYOFF}(\mathbf{x}, \mathbf{y})$$

$$(\mathbf{x}^*, \mathbf{y}^*) = \text{let } \mathbf{x}^* \stackrel{\Delta}{=} \text{GRADIENTDESCENT} \\ \quad \quad \quad ((\lambda \mathbf{x} \text{ GRADIENTDESCENT} ((\lambda \mathbf{y} (-\text{PAYOFF}(\mathbf{x}, \mathbf{y}))), \mathbf{y}_0, \epsilon)), \mathbf{x}_0, \epsilon) \\ \quad \text{in } (\mathbf{x}^*, (\text{GRADIENTDESCENT} ((\lambda \mathbf{y} (-\text{PAYOFF}(\mathbf{x}^*, \mathbf{y}))), \mathbf{y}_0, \epsilon)))$$

von Neumann, J. and Morgenstern, O. (1944). *Theory of Games and Economic Behavior*. Princeton University Press, Princeton, NJ.

Carl Gauss  
|  
Christoph Gudermann  
|  
Karl Weierstrass  
|  
Hermann Schwarz  
|  
Leopold Fejér  
|  
John von Neumann

# Function Inversion

$$f^{-1} y \triangleq \text{ROOT } ((\lambda x \ |f x| - y|), x_0, \epsilon)$$

# Neural Nets

$$\text{NEURON}(\mathbf{w}, \mathbf{x}) \triangleq \text{SIGMOID}(\mathbf{w} \cdot \mathbf{x})$$

$$\text{NEURALNET}([\mathbf{w}''; \mathbf{w}'_1; \dots; \mathbf{w}'_m], \mathbf{x}) \triangleq \\ \text{NEURON}(\mathbf{w}'', [\text{NEURON}(\mathbf{w}'_1, \mathbf{x}); \dots; \text{NEURON}(\mathbf{w}'_m, \mathbf{x})])$$

$$\text{ERROR } \mathbf{w} \triangleq \\ \|[y_1; \dots; y_n] - [\text{NEURALNET}(\mathbf{w}, \mathbf{x}_1); \dots; \text{NEURALNET}(\mathbf{w}, \mathbf{x}_n)]\|$$

$$\text{GRADIENTDESCENT}(\text{ERROR}, \mathbf{w}_0, \epsilon)$$

Rumelhart, D. E., Hinton, G. E., and Williams, R. J. (1986). Learning representations by back-propagating errors. *Nature*, **323**:533–6.

# Supervised Machine Learning

## Function Approximation

$$\text{ERROR } \theta \triangleq \|[y_1; \dots; y_n] - [f(\theta, \mathbf{x}_1); \dots; f(\theta, \mathbf{x}_n)]\|$$

GRADIENTDESCENT (ERROR,  $\theta_0, \epsilon$ )

# Maximum Likelihood Estimation

$$\text{GRADIENTDESCENT} \left( \left( \left( \lambda \theta \left( - \prod_{\mathbf{x} \in \mathcal{X}} P(\mathbf{x}|\theta) \right) \right) \right), \theta_0, \epsilon \right)$$

Fisher, R. A. (1922). On the mathematical foundations of theoretical statistics. *Philos. Trans. Roy. Soc. London Ser. A*, **222**:309–68.

# Engineering Design

```

PERFORMANCEOF SPLINECONTROLPOINTS  $\triangleq$ 
  let WING  $\triangleq$  SPLINETOSURFACE SPLINECONTROLPOINTS;
      AIRFLOW  $\triangleq$  PDESOLVER (WING, NAVIERSTOKES);
      LIFT, DRAG  $\triangleq$  SURFACEINTEGRAL (WING, AIRFLOW, FORCE);
      PERFORMANCE  $\triangleq$  DESIGNMETRIC (LIFT, DRAG, (WEIGHT WING))
in PERFORMANCE
  
```

```

GRADIENTDESCENT (PERFORMANCEOF, SPLINECONTROLPOINTS0,  $\epsilon$ )
  
```

# Work in Progress

Lambda: the ultimate calculus

Part of larger project to combine  $\lambda$ -calculus and AD.

# Work in Progress

Lambda: the ultimate calculus

Part of larger project to combine  $\lambda$ -calculus and AD.

- include both forward-mode and reverse-mode operators

# Work in Progress

Lambda: the ultimate calculus

Part of larger project to combine  $\lambda$ -calculus and AD.

- include both forward-mode and reverse-mode operators
- get numerically correct answers
- via correct transformation (time/space complexity)

# Work in Progress

## Lambda: the ultimate calculus

Part of larger project to combine  $\lambda$ -calculus and AD.

- include both forward-mode and reverse-mode operators
- get numerically correct answers
- via correct transformation (time/space complexity)
- make  $\vec{\mathcal{J}}$  and  $\overleftarrow{\mathcal{J}}$  first class
- make them apply to *all* (differentiable) programs
- handle nested application (self-application of  $\vec{\mathcal{J}}$  and  $\overleftarrow{\mathcal{J}}$ )

# Work in Progress

## Lambda: the ultimate calculus

Part of larger project to combine  $\lambda$ -calculus and AD.

- include both forward-mode and reverse-mode operators
- get numerically correct answers
- via correct transformation (time/space complexity)
- make  $\vec{\mathcal{J}}$  and  $\overleftarrow{\mathcal{J}}$  first class
- make them apply to *all* (differentiable) programs
- handle nested application (self-application of  $\vec{\mathcal{J}}$  and  $\overleftarrow{\mathcal{J}}$ )
- good puns

# Work in Progress

## Lambda: the ultimate calculus

Part of larger project to combine  $\lambda$ -calculus and AD.

- include both forward-mode and reverse-mode operators
- get numerically correct answers
- via correct transformation (time/space complexity)
- make  $\vec{\mathcal{J}}$  and  $\overleftarrow{\mathcal{J}}$  first class
- make them apply to *all* (differentiable) programs
- handle nested application (self-application of  $\vec{\mathcal{J}}$  and  $\overleftarrow{\mathcal{J}}$ )
- good puns
  - mathematical: the  $\lambda\nabla$ -calculus

# Work in Progress

## Lambda: the ultimate calculus

Part of larger project to combine  $\lambda$ -calculus and AD.

- include both forward-mode and reverse-mode operators
- get numerically correct answers
- via correct transformation (time/space complexity)
- make  $\vec{\mathcal{J}}$  and  $\overleftarrow{\mathcal{J}}$  first class
- make them apply to *all* (differentiable) programs
- handle nested application (self-application of  $\vec{\mathcal{J}}$  and  $\overleftarrow{\mathcal{J}}$ )
- good puns
  - mathematical: the  $\lambda\nabla$ -calculus
  - programming language: Voice for FP + AD

# Work in Progress

## Lambda: the ultimate calculus

Part of larger project to combine  $\lambda$ -calculus and AD.

- include both forward-mode and reverse-mode operators
- get numerically correct answers
- via correct transformation (time/space complexity)
- make  $\vec{\mathcal{J}}$  and  $\overleftarrow{\mathcal{J}}$  first class
- make them apply to *all* (differentiable) programs
- handle nested application (self-application of  $\vec{\mathcal{J}}$  and  $\overleftarrow{\mathcal{J}}$ )
- good puns
  - mathematical: the  $\lambda\nabla$ -calculus
  - programming language: VLAD

# Work in Progress

## Lambda: the ultimate calculus

Part of larger project to combine  $\lambda$ -calculus and AD.

- include both forward-mode and reverse-mode operators
- get numerically correct answers
- via correct transformation (time/space complexity)
- make  $\vec{\mathcal{J}}$  and  $\overleftarrow{\mathcal{J}}$  first class
- make them apply to *all* (differentiable) programs
- handle nested application (self-application of  $\vec{\mathcal{J}}$  and  $\overleftarrow{\mathcal{J}}$ )
- good puns
  - mathematical: the  $\lambda\nabla$ -calculus
  - programming language: VLAD
  - implementation: high-performance SCHEME + AD

# Work in Progress

## Lambda: the ultimate calculus

Part of larger project to combine  $\lambda$ -calculus and AD.

- include both forward-mode and reverse-mode operators
- get numerically correct answers
- via correct transformation (time/space complexity)
- make  $\vec{\mathcal{J}}$  and  $\overleftarrow{\mathcal{J}}$  first class
- make them apply to *all* (differentiable) programs
- handle nested application (self-application of  $\vec{\mathcal{J}}$  and  $\overleftarrow{\mathcal{J}}$ )
- good puns
  - mathematical: the  $\lambda\nabla$ -calculus
  - programming language: VLAD
  - implementation: STALIN + AD

# Work in Progress

## Lambda: the ultimate calculus

Part of larger project to combine  $\lambda$ -calculus and AD.

- include both forward-mode and reverse-mode operators
- get numerically correct answers
- via correct transformation (time/space complexity)
- make  $\vec{\mathcal{J}}$  and  $\overleftarrow{\mathcal{J}}$  first class
- make them apply to *all* (differentiable) programs
- handle nested application (self-application of  $\vec{\mathcal{J}}$  and  $\overleftarrow{\mathcal{J}}$ )
- good puns
  - mathematical: the  $\lambda\nabla$ -calculus
  - programming language: VLAD
  - implementation: STALIN + gradients

# Work in Progress

## Lambda: the ultimate calculus

Part of larger project to combine  $\lambda$ -calculus and AD.

- include both forward-mode and reverse-mode operators
- get numerically correct answers
- via correct transformation (time/space complexity)
- make  $\vec{\mathcal{J}}$  and  $\overleftarrow{\mathcal{J}}$  first class
- make them apply to *all* (differentiable) programs
- handle nested application (self-application of  $\vec{\mathcal{J}}$  and  $\overleftarrow{\mathcal{J}}$ )
- good puns
  - mathematical: the  $\lambda\nabla$ -calculus
  - programming language: VLAD
  - implementation: STALIN +  $\nabla$

# Work in Progress

## Lambda: the ultimate calculus

Part of larger project to combine  $\lambda$ -calculus and AD.

- include both forward-mode and reverse-mode operators
- get numerically correct answers
- via correct transformation (time/space complexity)
- make  $\vec{\mathcal{J}}$  and  $\overleftarrow{\mathcal{J}}$  first class
- make them apply to *all* (differentiable) programs
- handle nested application (self-application of  $\vec{\mathcal{J}}$  and  $\overleftarrow{\mathcal{J}}$ )
- good puns
  - mathematical: the  $\lambda\nabla$ -calculus
  - programming language: VLAD
  - implementation: STALIN $\nabla$

# Work in Progress

Lambda: the ultimate calculus

Part of larger project to combine  $\lambda$ -calculus and AD.

- include both forward-mode and reverse-mode operators
- get numerically correct answers
- via correct transformation (time/space complexity)
- make  $\vec{\mathcal{J}}$  and  $\overleftarrow{\mathcal{J}}$  first class
- make them apply to *all* (differentiable) programs
- handle nested application (self-application of  $\vec{\mathcal{J}}$  and  $\overleftarrow{\mathcal{J}}$ )
- good puns
  - mathematical: the  $\lambda\nabla$ -calculus
  - programming language: VLAD
  - implementation: STALIN $\nabla$
- manuscripts and code:

<http://www-bcl.cs.nuim.ie/~qobi/stalingrad/>