

Optimization of Sparse Generalized Jacobian Chain Products

Tamme Claus[†] and Uwe Naumann[‡]

[†] (tamme.claus@rwth-aachen.de)

[‡] supervisor, STCE (naumann@stce.rwth-aachen.de)

July 16, 2020

CES-Seminar, Master Computational Engineering Science,
RWTH Aachen University

What is Algorithmic Differentiation?

Sparse Generalized Jacobian Chain Product

Some Results

Motivation: Possible Applications of Algorithmic Differentiation

Sensitivity Analysis
Uncertainty Quantification
Derivatives

Machine Learning Gradient Descent
Optimization Problem
Inverse Problems Parameter Estimation

What is Algorithmic Differentiation?

Ways to calculate Derivatives

$$\mathbf{y} = F(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}^m$$

- Symbolic differentiation / Differentiation by hand

$$F'_{i,j}(\mathbf{x}) = \lim_{h \rightarrow 0} \frac{F_i(\mathbf{x} + h\mathbf{e}_j) - F_i(\mathbf{x})}{h}$$

→ exact, but **tedious** for arbitrary (complex) F

- Finite differences

$$F'_{i,j}(\mathbf{x}; h) \approx \frac{F_i(\mathbf{x} + h\mathbf{e}_j) - F_i(\mathbf{x})}{h}$$

→ **approximation**, choice of h

- Algorithmic/Computational Differentiation

$$\dot{v}_j = \sum_{k \prec j} \frac{\partial \varphi_j}{\partial v_k} \dot{v}_k \quad \bar{v}_k = \sum_{j \succ k} \frac{\partial \varphi_j}{\partial v_k} \bar{v}_j$$

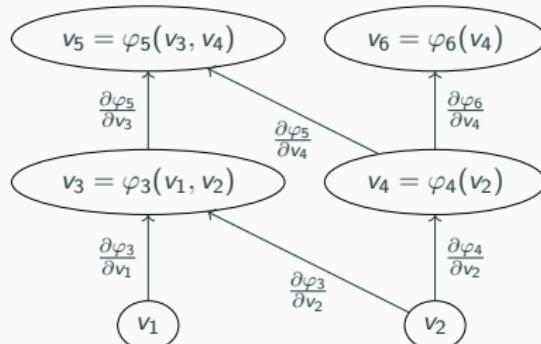
→ what? but **exact**, and **automatic**

Algorithmic/Computational Differentiation

- subdivide F in **elemental functions**

$$v_j = \varphi_j(\{v_k\}_{k \prec j}) \quad v_i \in \mathbb{R}$$

- define **tangents** \dot{v} and **adjoints** \bar{v}
- use the **chain rule**



$$\dot{v}_j = \sum_{k \prec j} \frac{\partial \varphi_j}{\partial v_k} \dot{v}_k$$

$$\bar{v}_k = \sum_{j \succ k} \frac{\partial \varphi_j}{\partial v_k} \bar{v}_j$$

$k \prec j$: all v_k which v_j depends on (predecessors)

$j \succ k$: all v_j which depend on v_k (successors)

→ still the goal: **Jacobian accumulation** but efficient!?

AD Mission Planning

Reducing fmas (fused multiply add-operations)

- occur frequently in AD
 - native support by modern processors
-

the time spent on

finding the optimum + optimal accumulation < naive accumulation

graph reduction methods: vertex/edge/face elimination

fine level of granularity, but intractable for large functions/computer code

here: function compositions

coarser level of granularity, restriction of the underlying problem,
focus on solvability, still enormous fma reductions (see later)

Sparse Generalized Jacobian Chain Product

Problem Definition

Sparse Generalized Jacobian Chain Product

$$F = F_q \circ F_{q-1} \circ \dots \circ F_1$$

$$F' = F'_q \cdot F'_{q-1} \cdot \dots \cdot F'_1$$

Assume that for each $z_i = F_i(z_{i-1}) : \mathbb{R}^{n_i} \rightarrow \mathbb{R}^{m_i}$

- tangent $\dot{F}_i \cdot \dot{Z}_{i-1}$ implementations,
- adjoint $\bar{Z}_i \cdot \bar{F}_i$ implementations, and
- sparsity pattern $S_i \in \{0, 1\}^{m_i \times n_i}$

are available.

What is the optimal way (minimum number of fmas) to accumulate F' ? (without cost reductions "inside" F_i)

Individual Costs: Naive Accumulation

Costs to accumulate $F'_i(\text{naive})$

- use tangent/adjoint implementation, "seed" unit vectors
- cost tangent: $n_i|E_i|$ ($|E_i|$ number of edges in the DAG of F_i)
- cost adjoint: $m_i|E_i|$

Individual Costs: Jacobian Compression (unidirectional)

Consider the sparsity pattern S_i of a factor F_i

$$S_i = \begin{pmatrix} 1 & 0 & 1 \\ 0 & 1 & 0 \end{pmatrix}$$

Tangent Compression

$$S_i \underbrace{\begin{pmatrix} 1 & 0 \\ 1 & 0 \\ 0 & 1 \end{pmatrix}}_{T_i^{(s)}(n_i \times t_i^{(s)})} = \begin{pmatrix} 1 & 1 \\ 1 & 0 \end{pmatrix}$$

→ graph coloring (heuristics)

Adjoint Compression

$$\underbrace{\begin{pmatrix} 1 & 1 \end{pmatrix}}_{A_i^{(s)}(a_i^{(s)} \times m_i)} S_i = \begin{pmatrix} 1 & 1 & 1 \end{pmatrix}$$

Costs

$$\text{fma}_{i,i} = |E_i| \cdot \min\left\{\min_{s \in \mathcal{S}_i^a}\{a_i^{(s)}\}, \min_{s \in \mathcal{S}_i^t}\{t_i^{(s)}\}\right\}$$

Individual Costs: Sparse Matrix Multiplication

Matrix Multiplication

Consider the sparsity pattern $A \in \{0, 1\}^{m_2 \times n_2}$ and $B \in \{0, 1\}^{m_1 \times n_1}$ of two matrices ($n_2 = m_1$)

$$\text{fma} = \sum_{i=1}^{m_2} \sum_{j=1}^{n_1} \sum_{k=1}^{n_2} A_{ik} B_{kj}$$

→ we only count multiplication of two non-zeros (assuming a sparse implementation)

Matrix Recompression

Multiplication of a matrix with a seed matrix: no multiplications, only additions:

→ count one *fma* per addition

Accumulation Options

Consider some subchain $F'_{j,i} \in \mathbb{R}^{m_j \times n_i}$

$$F'_{j,i} = \underbrace{F'_j \cdot \dots \cdot F'_{k+1}}_{F'_{j,k+1}} \cdot \underbrace{F'_k \cdot \dots \cdot F'_i}_{F'_{k,i}}$$

- **preaccumulate** both $F'_{j,k+1} \cdot F'_{k,i}$
(accumulation of both Jacobians + matrix multiplication)
- **tangent** mode $\dot{F}_{j,i} \cdot T_{j,i}^{(s)}$ or $\dot{F}_{j,k+1} \cdot (F'_{k,i} \cdot T_{j,i}^{(s)})$
(accumulation of Jacobian + matrix recompression + cost of tangent)
- **adjoint** mode $A_{j,i}^{(s)} \cdot \bar{F}_{j,i}$ or $(A_{j,i}^{(s)} \cdot F'_{j,k+1}) \cdot \bar{F}_{k,i}$
(accumulation of Jacobian + matrix recompression + cost of adjoint)

Dynamic Programming

overlapping subproblems, optimal substructure

→ optimal solution of smaller Jacobian chains may be reused

Dynamic Programming Recurrence

- preaccumulate both $F'_{j,k+1} \cdot F'_{k,i}$
- tangent mode $\dot{F}_{j,i} \cdot T_{j,i}^{(s)}$ or $\dot{F}_{j,k+1} \cdot (F'_{k,i} \cdot T_{j,i}^{(s)})$
- adjoint mode $A_{j,i}^{(s)} \cdot \bar{F}_{j,i}$ or $(A_{j,i}^{(s)} \cdot F'_{j,k+1}) \cdot \bar{F}_{k,i}$

Recurrence

$$\text{fma}_{j,i} = \begin{cases} |E_j| \cdot \min\left\{\min_{s \in S_{j,j}^a}\{a_j^{(s)}\}, \min_{s \in S_{j,j}^t}\{t_j^{(s)}\}\right\} & i = j \\ \min\left\{\min_{i \leq k < j}\left\{\min\left\{\begin{array}{l} \text{fma}_{j,k+1} + \text{fma}_{k,i} + \text{fma}_{j,k,i}, \\ \text{fma}_{j,k+1} + \min_{s \in S_{j,i}^a}\{\text{fma}_{j,k+1,i}^{(s)} + a_{j,i}^{(s)} \sum_{\nu=i}^k |E_\nu|\}, \\ \text{fma}_{k,i} + \min_{s \in S_{j,i}^t}\{\text{fma}_{j,k,i}^{(s)} + t_{j,i}^{(s)} \sum_{\nu=k+1}^j |E_\nu|\} \end{array}\right\}\right\}\right\} & j > i \end{cases}$$

Restrictions

the time spent on

finding the optimum + optimal accumulation < naive accumulation

Balance between

- rigorous minimal fma
- time spent optimizing

→ restrict to one Jacobian compression (e.g. the one found by a heuristic)

→ do not consider Jacobian compression
(see results later)

Example

$$F_2^t \in \mathbb{R}^{4 \times 3} \quad |E_2| = 12$$

$$\begin{vmatrix} 0 & 1 & 0 \\ 1 & 0 & 1 \\ 1 & 1 & 1 \\ 0 & 1 & 0 \\ 0 & b & c \end{vmatrix}^a$$

$$\text{fma}_5 = |E_2| \min\{3, 3\} = 36$$

$$F_4^t \in \mathbb{R}^{3 \times 3} \quad |E_4| = 46$$

$$\begin{vmatrix} 1 & 1 & 1 \\ 1 & 0 & 1 \\ 1 & b & c \end{vmatrix}^a$$

$$\text{fma}_4 = |E_4| \min\{3, 3\} = 138$$

$$F_2^t \in \mathbb{R}^{3 \times 5} \quad |E_2| = 28$$

$$\begin{vmatrix} 1 & 1 & 1 & 1 & 0 \\ 1 & 1 & 0 & 1 & 0 \\ 1 & 0 & 1 & 0 & 1 \\ 1 & a & b & c & d \\ a & b & c & d & b \end{vmatrix}^a$$

$$\text{fma}_3 = |E_2| \min\{3, 4\} = 84$$

$$F_2^t \in \mathbb{R}^{3 \times 2} \quad |E_2| = 15$$

$$\begin{vmatrix} 0 & 1 \\ 1 & 0 \\ 0 & 1 \\ 0 & 1 \\ a & a \end{vmatrix}^a$$

$$\begin{vmatrix} 1 & 0 & a \\ 0 & 1 & b \\ 0 & a & b \end{vmatrix}^a$$

$$\text{fma}_1 = |E_1| \min\{1, 2\} = 5$$

$$F_{3,4}^t \in \mathbb{R}^{4 \times 3}$$

$$\begin{vmatrix} 1 & 0 & 1 & a \\ 1 & 1 & 1 & b \\ 1 & 1 & 1 & c \\ 1 & 1 & 1 & d \\ a & b & c & e \end{vmatrix}^a$$

$$\text{fma}_{4,4} = \min \left\{ \begin{array}{l} (|E_3| + |E_4|) \min\{4, 3\} \\ \text{fma}_5 + \text{fma}_4 + 16 \\ \text{fma}_5 + 0 + |E_4| \\ \text{fma}_4 + 0 + |E_3|^2 \end{array} \right\} = 174$$

$$F_{3,3}^t \in \mathbb{R}^{3 \times 5}$$

$$\begin{vmatrix} 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 & 0 \\ a & b & c & d & e \end{vmatrix}^a$$

$$\text{fma}_{4,3} = \min \left\{ \begin{array}{l} (|E_3| + |E_4|) \min\{3, 5\} \\ \text{fma}_5 + \text{fma}_4 + 24 \\ \text{fma}_5 + 0 + |E_3|^3 \\ \text{fma}_3 + 0 + |E_5|^5 \end{array} \right\} = 222$$

$$F_{3,2}^t \in \mathbb{R}^{3 \times 2}$$

$$\begin{vmatrix} 1 & 1 \\ 1 & 0 \\ 0 & 1 \\ 0 & 1 \\ a & b \end{vmatrix}^a$$

$$\text{fma}_{3,2} = \min \left\{ \begin{array}{l} (|E_3| + |E_2|) \ min\{3, 2\} \\ \text{fma}_5 + \text{fma}_4 + 10 \\ \text{fma}_5 + 0 + |E_2|^3 \\ \text{fma}_2 + 0 + |E_2|^2 \end{array} \right\} = 71$$

$$F_{2,1}^t \in \mathbb{R}^{3 \times 3}$$

$$\begin{vmatrix} 0 & 1 & 0 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \\ 0 & 1 & 0 \\ 0 & 1 & 1 \\ a & a & b \end{vmatrix}^a$$

$$\text{fma}_{2,1} = \min \left\{ \begin{array}{l} (|E_2| + |E_1|) \ min\{3, 2\} \\ \text{fma}_5 + \text{fma}_4 + 8 \\ \text{fma}_2 + 0 + |E_1|^3 \\ \text{fma}_1 + 0 + |E_2|^2 \end{array} \right\} = 28$$

$$F_{3,3}^t \in \mathbb{R}^{4 \times 5}$$

$$\begin{vmatrix} 1 \\ 1 \end{vmatrix}$$

$$\text{fma}_{3,3} = \min \left\{ \begin{array}{l} (|E_3| + |E_4| + |E_5|) \ min\{4, 5\} \\ \min\{\text{fma}_5 + \text{fma}_4 + 34, \text{fma}_5 + \text{fma}_{4,4} + 33\} \\ \min\{\text{fma}_4 + 0 + (|E_4| + |E_5|) 1, \text{fma}_{4,4} + 0 + |E_4|^4\} \\ \min\{\text{fma}_4 + 0 + (|E_5| + |E_4|) 5, \text{fma}_{4,3} + 0 + |E_5|^5\} \end{array} \right\} = 282$$

$$F_{4,2}^t \in \mathbb{R}^{3 \times 2}$$

$$\begin{vmatrix} 1 \\ 1 \end{vmatrix}$$

$$\text{fma}_{4,2} = \min \left\{ \begin{array}{l} (|E_4| + |E_5|) \ min\{3, 2\} \\ \min\{\text{fma}_4 + \text{fma}_3 + 14, \text{fma}_4 + \text{fma}_{3,2} + 14\} \\ \min\{\text{fma}_4 + 0 + (|E_4| + |E_5|) 3, \text{fma}_{4,3} + 0 + |E_4|^3\} \\ \min\{\text{fma}_2 + 0 + (|E_4| + |E_5|) 2, \text{fma}_{3,2} + 0 + |E_4|^2\} \end{array} \right\} = 163$$

$$F_{3,1}^t \in \mathbb{R}^{3 \times 3}$$

$$\begin{vmatrix} 1 \\ 1 \end{vmatrix}$$

$$\text{fma}_{3,1} = \min \left\{ \begin{array}{l} (|E_3| + |E_2| + |E_1|) \ min\{3, 3\} \\ \min\{\text{fma}_{3,2} + \text{fma}_2 + 9, \text{fma}_3 + \text{fma}_{2,1} + 15\} \\ \min\{\text{fma}_4 + 0 + (|E_2| + |E_1|) 3, \text{fma}_{4,2} + 0 + |E_1|^3\} \\ \min\{\text{fma}_4 + 0 + (|E_3| + |E_2|) 3, \text{fma}_{3,2} + 0 + |E_3|^3\} \end{array} \right\} = 85$$

$$F_{3,2}^t \in \mathbb{R}^{4 \times 2}$$

$$\begin{vmatrix} 1 \\ 1 \end{vmatrix}$$

$$\text{fma}_{5,2} = \min \left\{ \begin{array}{l} (|E_5| + |E_4| + |E_3|) \ min\{4, 2\} \\ \min\{\text{fma}_5 + \text{fma}_4 + 20, \text{fma}_5 + \text{fma}_{3,2} + 20, \text{fma}_5 + \text{fma}_{4,2} + 14\} \\ \min\{\text{fma}_2 + 0 + (|E_4| + |E_3| + |E_2|) 4, \text{fma}_{4,4} + 0 + (|E_3| + |E_2|) 4, \text{fma}_{5,2} + 0 + |E_2|^4\} \\ \min\{\text{fma}_2 + 0 + (|E_5| + |E_4| + |E_3|) 2, \text{fma}_{3,2} + 0 + (|E_5| + |E_4|) 2, \text{fma}_{4,2} + 0 + |E_3|^2\} \end{array} \right\} = 187$$

$$F_{5,2}^t \in \mathbb{R}^{4 \times 3}$$

$$\begin{vmatrix} 1 \\ 1 \end{vmatrix}$$

$$\text{fma}_{5,1} = \min \left\{ \begin{array}{l} (|E_5| + |E_4| + |E_3| + |E_2| + |E_1|) \ min\{4, 3\} \\ \min\{\text{fma}_{5,2} + \text{fma}_4 + 12, \text{fma}_{5,3} + 32, \text{fma}_{4,4} + \text{fma}_{3,1} + 30, \text{fma}_5 + \text{fma}_{4,1} + 21\} \\ \min\{\text{fma}_5 + 0 + (|E_4| + |E_3| + |E_2| + |E_1|) 4, \text{fma}_{5,4} + 0 + (|E_3| + |E_2| + |E_1|) 4, \text{fma}_{5,3} + 0 + (|E_2| + |E_1|) 4, \text{fma}_{5,2} + 0 + |E_1|^4\} \\ \min\{\text{fma}_1 + 0 + (|E_5| + |E_4| + |E_3| + |E_2| + |E_1|) 3, \text{fma}_{2,3} + 0 + (|E_5| + |E_4| + |E_3| + |E_2|) 3, \text{fma}_{3,1} + 0 + (|E_5| + |E_4| + |E_3| + |E_2|) 3, \text{fma}_{4,1} + 0 + |E_3|^3\} \end{array} \right\} = 204$$

$$F_{4,1}^t \in \mathbb{R}^{3 \times 3}$$

$$\begin{vmatrix} 1 \\ 1 \end{vmatrix}$$

$$\text{fma}_{4,1} = \min \left\{ \begin{array}{l} (|E_4| + |E_3| + |E_2| + |E_1|) \ min\{3, 3\} \\ \min\{\text{fma}_4 + \text{fma}_3 + 9, \text{fma}_{4,3} + \text{fma}_{2,1} + 22, \text{fma}_4 + \text{fma}_{3,1} + 21\} \\ \min\{\text{fma}_5 + 0 + (|E_4| + |E_3| + |E_2| + |E_1|) 3, \text{fma}_{5,3} + 0 + (|E_3| + |E_2| + |E_1|) 3, \text{fma}_{4,2} + 0 + |E_1|^3\} \\ \min\{\text{fma}_5 + 0 + (|E_3| + |E_2| + |E_1|) 3, \text{fma}_{4,3} + 0 + (|E_2| + |E_1|) 3, \text{fma}_{3,2} + 0 + (|E_2| + |E_1|) 3, \text{fma}_{2,1} + 0 + (|E_2| + |E_1|) 3\} \end{array} \right\} = 177$$

Some Results

Case Study

Dense

length q	max_mn	Tangent	Adjoint	Preaccumulation	Optimum
10	10	296	2072	779	296
50	50	839925	604746	740518	135228
100	100	19160886	21611697	9034346	231143
150	150	56863838	47527984	53425505	852466
200	200	34648499	381133489	174604509	2041326

Sparse - no compression

10	10	296	2072	681	296
50	50	839925	604746	620148	134838
100	100	19160886	21611697	8808143	231143
150	150	56863838	47527984	52552009	852466
200	200	34648499	381133489	172585171	2041326

Sparse - compression

10	10	296	2072	345	296
50	50	839925	604746	36767	36132
100	100	19160886	21611697	237339	231143
150	150	56863838	47527984	867098	852466
200	200	34648499	381133489	2067283	2041326

Case Study

Dense

length q	max_mn	Tangent	Adjoint	Preaccumulation	Optimum
10	10	296	2K	779	296
50	50	840K	605K	741K	135K
100	100	19M	22M	9M	231K
150	150	57M	48M	53M	852K
200	200	35M	381M	175M	2M

Sparse - no compression

10	10	296	2K	681	296
50	50	840K	605K	620K	135K
100	100	19M	22M	9M	231K
150	150	57M	48M	53M	852K
200	200	35M	381M	173M	2M

Sparse - compression

10	10	296	2K	345	296
50	50	840K	605K	37K	36K
100	100	19M	22M	237K	231K
150	150	57M	48M	867K	852K
200	200	35M	381M	2M	2M

Outlook

Conclusion

- enormous fma improvements possible compared to naive AD
- shows potential and necessity of *AD Mission Planning*
- exploitation of sparsity: effectiveness is problem dependent

Future work

- adaptive optimization finding balance between degrees of freedom and optimization time
- memory constraints for adjoints
- extension to different function structure (more than function compositions)

References

- Gebremedhin, A., Manne, F., & Pothen, A. (2005). What color is your jacobian? graph coloring for computing derivatives. *SIAM review*, 47, 629–705. doi:10.1137/S0036144504444711
- Giering, R., & Kaminski, T. (2006). Automatic sparsity detection implemented as a source-to-source transformation. In V. N. Alexandrov, G. D. van Albada, P. M. A. Sloot, & J. Dongarra (Eds.), *Computational science – iccs 2006* (pp. 591–598). Berlin, Heidelberg: Springer Berlin Heidelberg.
- Goedecker, S., & Hoisie, A. (2001). *Performance optimization of numerically intensive codes*. doi:10.1137/1.9780898718218. eprint: <https://epubs.siam.org/doi/pdf/10.1137/1.9780898718218>
- Griewank, A. (2003). A mathematical view of automatic differentiation. *Acta Numerica*, 12, 321–398. doi:10.1017/S0962492902000132
- Griewank, A., & Mitev, C. (2002). Detecting jacobian sparsity patterns by bayesian probing. *Math. Program.* 93, 1–25. doi:10.1007/s101070100281
- Griewank, A., & Naumann, U. (2002). Accumulating jacobians as chained sparse matrix products. *Math. Prog.* 3, 2003.
- Naumann, U. (2004). Optimal accumulation of jacobian matrices by elimination methods on the dual computational graph. *Math. Program.* 99, 399–421. doi:10.1007/s10107-003-0456-9
- Naumann, U. (2008). Optimal jacobian accumulation is np-complete. *Mathematical Programming*, 112, 427–441. doi:10.1007/s10107-006-0042-z
- Naumann, U. (2011). *The art of differentiating computer programs*. doi:10.1137/1.9781611972078
- Naumann, U. (2020a). On sparse matrix chain products. (pp. 118–127). doi:10.1137/1.9781611976229.12
- Naumann, U. (2020b). Optimization of generalized jacobian chain products without memory constraints. arXiv: 2003.05755 [math.NA]
- Pieterse, V., Kourie, D., Cleophas, L., & Watson, B. (2010). Performance of c++ bit-vector implementations. (pp. 242–250). doi:10.1145/1899503.1899530