

590



Derivative Code by Overloading in C++

... with a clear focus on adjoints and building on Laurent's talk

Uwe Naumann

RWTH Aachen University, Germany naumann@stce.rwth-aachen.de

▲□▶ ▲□



Basic Adjoint AD

First Order Implementation Case Study

Beyond Basic Adjoint AD

Checkpointing Preaccumulation Higher-Level Adjoint Elementals Debugging Adjoints Further Aspects

Conclusion

Recent Applications / Publications

Miscellaneous

Inside 1st-Order Algorithmic Adjoints Beyond 1st-Order Adjoints Further Details

Outline



Basic Adjoint AD

First Order Implementation Case Study

Beyond Basic Adjoint AD

Checkpointing Preaccumulation Higher-Level Adjoint Elementals Debugging Adjoints Further Aspects

Conclusion

Recent Applications / Publications

Miscellaneous

Inside 1st-Order Algorithmic Adjoints Beyond 1st-Order Adjoints Further Details

Algorithmic Differentiation (AD)



First-Order Adjoints

Let the **C++ program** $\mathbf{y} := F(\mathbf{x}) : \mathbb{R}^n \to \mathbb{R}^m$ be continuously differentiable at all points of interest with Jacobian $\nabla F(\mathbf{x}) \equiv \frac{dF}{d\mathbf{x}}(\mathbf{x}) \in \mathbb{R}^{m \times n}$. Typically, Fimplements the numerical approximation of a function $\mathbf{y}(\mathbf{x})$ implicitly defined via $R(\mathbf{y}(\mathbf{x})) = 0$.

Overloading and template metaprogramming techniques, e.g, in $dco/c++,^1$ yield the **adjoint program**

 $\mathbf{x}_{(1)} \mathrel{+}= \nabla F(\mathbf{x})^T \cdot \mathbf{y}_{(1)}$

through use of a custom (e.g, generic <u>a</u>djoint <u>1</u>st-order <u>s</u>calar) data type std::vector<dco::ga1s<T>::type> x(n),y(m);

over base type T for recording a (e.g, global) tape

dco::ga1s<T>::tape_t* dco::ga1s<T>::global_tape;

and its subsequent interpretation.

¹www.nag.co.uk

 $\begin{array}{c} \text{Implementation} \\ \mathbf{y} = \underline{\mathbf{x}} := F(\mathbf{x}) = G(H(\mathbf{x}), \mathbf{x}) \quad \Rightarrow \quad \mathbf{x}_{(1)} := \nabla F(\mathbf{x})^T \cdot \mathbf{x}_{(1)} \end{array}$

Source Transformation

$$\begin{split} \mathbf{z} &:= H(\mathbf{x}) \\ \mathsf{PUSH}(\mathbf{x}) & // \text{ if } \mathsf{TBR}(\mathbf{x}) \\ \mathbf{x} &:= G(\mathbf{z}, \mathbf{x}) & // \text{ or move to end} \end{split}$$

$$\begin{aligned} & \mathsf{POP}(\mathbf{x}) \\ & \mathbf{z}_{(1)} \coloneqq \frac{dG}{d\mathbf{z}}^T \cdot \mathbf{x}_{(1)} \\ & \mathbf{x}_{(1)} \coloneqq \frac{\partial G}{\partial \mathbf{x}}^T \cdot \mathbf{x}_{(1)} \\ & \mathbf{x}_{(1)} \coloneqq \mathbf{x}_{(1)} + \frac{dH}{d\mathbf{x}}^T \cdot \mathbf{z}_{(1)} \end{aligned}$$

 $\textbf{Overloading}~(\rightarrow \mathsf{Tape}~[\mathsf{DAG}])$



$$\Leftarrow \quad \text{chain rule:} \quad \nabla F(\mathbf{x}) = \sum_{\text{path} \in \text{DAG}} \prod_{(i,j) \in \text{path}} d_{j,i}$$

Case Study Stochastic Differential Equation [5]



We are looking for the expected value $\mathbb{E}(x)$ of the solution x(p(t), T), $0 \le t \le T$, of the scalar stochastic initial value problem

 $dx = p(t) \cdot \sin(x(p(t), t) \cdot t)dt + p(t) \cdot \cos(x(p(t), t) \cdot t)dW; \ t \in [0, T]$

with Brownian Motion dW and for $x(p(0),0)=x^0.$ Forward finite differences in time with time step $0<\Delta t\ll T$ yield the Euler-Maruyama scheme

 $x^{i+1} := x^i + \Delta t \cdot p_i \cdot \sin(x^i \cdot i \cdot \Delta t) + \sqrt{\Delta t} \cdot p_i \cdot \cos(x^i \cdot i \cdot \Delta t) \cdot dW^i$

for $i = 0, \ldots, n-1$, target time $T = n \cdot \Delta t$, parameter vector $\mathbf{p} = (p_i) \in \mathbb{R}^l$, and with random numbers dW^i drawn from the standard normal distribution N(0, 1). The simulation of m Monte Carlo paths with precomputed random numbers $dW \in \mathbb{R}^{m \times n}$ is implemented as



Gradient with dco/c++: $\mathbf{p}_{(1)} := \frac{df}{d\mathbf{p}}(x, \mathbf{p})^T \cdot x_{(1)}$

```
#include "dco.hpp"
1
2
   template<typename AFLT_T, typename PFLT_T>
3
   void dxdp_a(AFLT_T& x_v, std::vector<AFLT_T>& p_v,
       const std::vector<std::vector<PFLT T>> &dW.
5
       std::vector<AFLT_T> &dxdp ) {
6
     typedef typename dco::gals<AFLT_T> DCO_M; // adjoint mode
7
     typedef typename DCO_M::type DCO_T; // adjoint type
8
     typedef typename DCO_M::tape_t DCO_TT; // tape type
9
     int n=p_v.size();
10
11
     DCO_T x=x_v;
      std::vector<DCO_T> p(n); dco::value(p)=p_v;
12
     DCO_M::global_tape=DCO_TT::create(); // create tape
13
     DCO_M::global_tape->register_variable(p); // record active input
14
     f(x,p,dW); // record active computation
15
     dco::derivative(x)=1: // seed
16
     DCO_M::global_tape->interpret_adjoint(); // interpret tape
17
     dxdp=dco::derivative(p); // harvest
18
     DCO_TT::remove(DCO_M::global_tape); // remove tape
19
   }
20
```

Case Study



| | Time (s) | Memory (MB) | rel. Time |
|---------------------------|----------|-------------|-----------|
| Primal | 1.2 | 79 | 1 |
| FFD | 380 | 80 | 317 |
| AAD by hand ² | 1.9 | 240 | 1.6 |
| pathwise AAD by hand | 2.2 | 80 | 1.8 |
| AAD by dco/c++ | 1.7 | 728 | 1.4 |
| pathwise AAD by $dco/c++$ | 1.9 | 86 | 1.6 |

FFD: Forward Finite Differences AAD: Adjoint Algorithmic Differentiation

²using basic adjoint code generation rules [1]



W.l.o.g, let the C++ program $y := f(\mathbf{x}) : \mathbb{R}^n \to \mathbb{R}$ be twice continuously differentiable with Hessian $\nabla^2 f(\mathbf{x}) \equiv \frac{d^2 f}{d\mathbf{x}^2}(\mathbf{x}) \in \mathbb{R}^{n \times n}$.

Nested first-order derivative types yield 2nd-order adjoints

 $\mathbf{x}_{(1)}^{(2)} \mathrel{+}= y_{(1)} \cdot \nabla^2 f(\mathbf{x}) \cdot \mathbf{x}^{(2)} \quad \left(\mathsf{Hessian at } O(n) \cdot \mathsf{Cost}(\mathbf{x}_{(1)}^{(2)}) \right),$

e.g, tangents of adjoints as

<dco::ga1s<dco::gt1s<T>::type>::type

Deeper nesting yields adjoints of arbitrary order, e.g, tangents of adjoints of adjoints as

dco::ga1s<dco::gt1s<T>::type>::type>::type

Outline



Basic Adjoint AD

First Order Implementation Case Study

Beyond Basic Adjoint AD

Checkpointing Preaccumulation Higher-Level Adjoint Elementals Debugging Adjoints Further Aspects

Conclusion

Recent Applications / Publications

Miscellaneous

Inside 1st-Order Algorithmic Adjoints Beyond 1st-Order Adjoints Further Details Beyond Basic Adjoint AD



"Playground", e.g, $\mathbf{x}=\mathbf{x}(\mathbf{p}):\mathbb{R}^n\to\mathbb{R}^m\Leftarrow R(\mathbf{x},\mathbf{p})=0$



POC $\hat{=}$ primal ops count (P)MR $\hat{=}$ (primal) memory requirement

Note various combinations and nesting.



The adjoint of a program $\mathbf{y} := F(\mathbf{x})$ computes

$$\boldsymbol{v}_{0(1)} = \underset{\in \mathbf{R}^n}{\mathbf{x}_{(1)}} \coloneqq \nabla F(\mathbf{x})^T \cdot \underset{\in \mathbf{R}^m}{\mathbf{y}_{(1)}} = \nabla \varphi_1^T \cdot (\dots (\nabla \varphi_q^T \cdot \boldsymbol{v}_{q(1)}) \dots)$$

assuming availability of adjoint elementals

$$\boldsymbol{v}_{i-1(1)} = \nabla \varphi_i (\boldsymbol{v}_{i-1})^T \cdot \boldsymbol{v}_{i(1)}$$

for $i = q, \ldots, 1$ (\rightarrow reversal of data flow).

The minimum requirement for adjoint AD (AAD) is the implementation of adjoint versions of the intrinsic operations (+, *, ...) and functions $(\sin, \exp, ...)$ of the given programming language.

Their naive combination yields algorithmic adjoint programs, which may turn out infeasible for various reasons. Hierarchies in granularity and mathematical semantics must be exploited in "real world" AAD.



Let $\varphi_{k(1)}$ not be implemented by basic AAD.

The corresponding gap in the tape needs to be filled by the adjoint program

$$\mathbf{x}_{(1)} = \boldsymbol{v}_{0(1)} \coloneqq \nabla \varphi_1^T \cdot \ldots \cdot \overline{\nabla \varphi_k^T(\boldsymbol{v}_{k-1}) \cdot \underbrace{(\nabla \varphi_{k+1}^T \cdot \ldots \cdot (\nabla \varphi_q^T \cdot \mathbf{y}_{(1)}) \ldots)}_{\boldsymbol{v}_{k(1)}}$$

to be filled by a custom version of $\varphi_{k(1)}$.

For example, checkpointing methods decrease the maximum tape size by storing v_{k-1} in the forward section followed by the evaluation of the primal φ_k and postponing the generation of the tape for $\varphi_{(1)_k}$ to the reverse section of $\varphi_{(1)}$.

Further examples include the implementation of symbolic adjoint elementals, preaccumulation, coupling with hand-written / generated adjoint source code and approximation of Jacobians of local black boxes by finite differences.

Adjoint Code Design Patterns



Recursive (e.g, Binomial) Checkpointing



Naumann, Paris, 24.01.2019

Preaccumulation Reducing OC and MR







Let the primal matrix-matrix product

 $Y \mathrel{\mathop:}= A \cdot X$

for $X \in \mathbb{R}^{n \times p}$, $Y \in \mathbb{R}^{m \times p}$, and $A \in \mathbb{R}^{m \times n}$ be implemented (naively) as a function Ax as follows

```
1 ...
2 for (int i=0;i<m;i++)
3 for (int j=0;j<n;j++)
4 for (int k=0;k<p;k++)
5 Y[i][k]+=A[i][j]*X[j][k];
6 ...
</pre>
```

Let $Y := A^T \cdot X$ be implemented as ATX and $Y := A \cdot X^T$ as AXT.



Application of the adjoint code generation rules yields the adjoint version AX_a

of ax implementing $A_{(1)} += Y_{(1)} \cdot X^T$, $X_{(1)} += A^T \cdot Y_{(1)}$.

Architecture-specific optimized implementations of AXT and ATX should be used for the computation of AX_a .

Higher-Level Adjoint Elementals $F(\mathbf{x}(\mathbf{p}), \mathbf{p}) = 0$



Systems of Linear Equations [3]

 $A \cdot \mathbf{x} = \mathbf{b}$

At the primal solution $\mathbf{x} = \mathbf{x}(A, \mathbf{b})$ we derive

$$\mathbf{b}_{(1)} := A^{-T} \cdot \mathbf{x}_{(1)}$$
$$A_{(1)} := -\mathbf{b}_{(1)} \cdot \mathbf{x}^{T}$$

Systems of Parameterized Nonlinear Equations [7]

 $F(\mathbf{x}, \mathbf{p}) = 0$

At the primal solution $\mathbf{x} = \mathbf{x}(\mathbf{p})$ the implicit function theorem yields

$$\mathbf{z}_{(1)} \coloneqq \frac{dF}{d\mathbf{x}}^{-T} \mathbf{x}_{(1)} ; \quad \mathbf{p}_{(1)} \coloneqq -\frac{\partial F}{\partial \mathbf{p}}^{T} \cdot \mathbf{z}_{(1)}$$

Naumann, Paris, 24.01.2019

Higher-Level Adjoint Elementals $\min_{\mathbf{x}} f(\mathbf{x}(\mathbf{p}), \mathbf{p})$



At the primal solution $\mathbf{x} = \mathbf{x}(\mathbf{p})$ the necessary first-order optimality condition

$$\frac{df}{d\mathbf{x}}(\mathbf{x},\mathbf{p}) = 0$$

yields

$$\mathbf{z}_{(1)} := \frac{d^2 f^{-T}}{d \mathbf{x}^2} \cdot \mathbf{x}_{(1)} \; ; \quad \mathbf{p}_{(1)} := -\frac{\partial^2 f}{\partial \mathbf{x} \partial \mathbf{p}}^T \cdot \mathbf{z}_{(1)} \; .$$

E.g, the NAG AD Library routine $e04gb_a1w$ implements symbolic <u>a</u>djoint of length <u>1</u> (scalar) using <u>w</u>orking precision for

- unconstrained minimum of a sum of squares
- combined Gauss–Newton and quasi-Newton algorithm
- first derivatives provided by user.

Higher-Level Adjoint Elementals min_x $f(\mathbf{x}(\mathbf{p}))$ s.t. $g(\mathbf{x}(\mathbf{p}), \mathbf{p}) = 0$



At the primal solution $(\mathbf{x},\lambda)=(\mathbf{x}(\mathbf{p}),\lambda(\mathbf{p}))$ the KKT condition

$$\frac{d}{d\mathbf{x}}\mathcal{L}(\mathbf{x},\mathbf{p}) \equiv \frac{d}{d\mathbf{x}}\left[f(\mathbf{x}) + \lambda^T \cdot g(\mathbf{x},\mathbf{p})\right] = 0$$

yields (with appropriate tensor product $\langle\cdot,\cdot\rangle$ [1])

$$\begin{aligned} \mathbf{z}_{(1)} &\coloneqq \left(\frac{d^2 f}{d\mathbf{x}^2} + \langle \lambda, \frac{d^2 g}{d\mathbf{x}^2} \rangle \right)^{-T} \cdot \mathbf{x}_{(1)} \\ \mathbf{p}_{(1)} &\coloneqq -\left(\frac{\partial^2 f}{\partial \mathbf{x} \partial \mathbf{p}} + \langle \lambda, \frac{\partial^2 g}{\partial \mathbf{x} \partial \mathbf{p}} \rangle \right)^T \cdot \mathbf{z}_{(1)} \;. \end{aligned}$$

Moreover

$$\frac{df}{d\mathbf{p}}^{T} = \frac{d\mathbf{x}}{d\mathbf{p}}^{T} \cdot \frac{df}{d\mathbf{x}}^{T} := \frac{\partial g}{\partial \mathbf{p}}^{T} \cdot \lambda$$

identifies the Lagrange multipliers λ as seed of the adjoint constraint for investigating sensitivities of the optimum with respect to **p**.

Naumann, Paris, 24.01.2019

Adjoint Code Design Patterns



Symbolic Adjoint $F(\mathbf{x}, \mathbf{p}) = 0$



[2]

Debugging Adjoints $\langle \mathbf{x}^{(1)}, \mathbf{x}_{(1)} \rangle = v^{(1)} \cdot v_{(1)}$



requires dco::gta1s<T>::type and primal type generic external adjoints $\varphi_{2(1)}$ and $\varphi_{3(1)}$.



- shared memory parallelism (e.g, MPI)
- distributed memory parallelism (e.g, OpenMP)
- vectorization (e.g, AVX)
- massively parallel acceleration (e.g, GPU)
- cross-language adjoints (e.g, LLVM)
- static/dynamic program analysis (e.g, MIP for CTR)
- globalization of adjoints (e.g, convex relaxations)

Outline



Basic Adjoint AD

First Order Implementation Case Study

Beyond Basic Adjoint AD

Checkpointing Preaccumulation Higher-Level Adjoint Elementals Debugging Adjoints Further Aspects

Conclusion

Recent Applications / Publications

Miscellaneous

Inside 1st-Order Algorithmic Adjoints Beyond 1st-Order Adjoints Further Details

Software and Tools for Computational Engineering

The quality of an adjoint AD solution / tool is defined by

- level of mathematical and software technological abstraction
- robustness wrt. language features of target code
- efficiency of adjoint propagation
- scalability on modern parallel computer architectures
- flexibility wrt. design scenarios
- sustainability wrt. dynamics in user requirements, personnel, hard- and software

Outline



Basic Adjoint AD

First Order Implementation Case Study

Beyond Basic Adjoint AD

Checkpointing Preaccumulation Higher-Level Adjoint Elementals Debugging Adjoints Further Aspects

Conclusion

Recent Applications / Publications

Miscellaneous

Inside 1st-Order Algorithmic Adjoints Beyond 1st-Order Adjoints Further Details

Naumann, Paris, 24.01.2019

Selected Recent Applications









[Towara, N. 2015]

[Towara 2018]



[Kopmann, Riehme, N. 2017]



[Gremse, Hoefter, N. 2016]



▶ $[1] \rightarrow$

- ▶ [2] U.N.: Adjoint Code Design Patterns, under review.
- ▶ [3] U.N.: A Note on Adjoint Linear Algebra, under review.
- [4] M. Towara, U.N.: SIMPLE Adjoint Message Passing, OMS 2018.
- [5] U.N., J. du Toit: Adjoint Algorithmic Differentiation Tool Support for Typical Numerical Patterns in Computational Finance, JCF 2018.
- [6] U.N., K. Leppkes: Low-Memory Algorithmic Adjoint Propagation, SIAM CSC, 2018.
- [7] U.N., K. Leppkes, J. Lotz, M. Towara: Algorithmic Differentiation of Numerical Methods: Tangent and Adjoint Solvers for Parameterized Systems of Nonlinear Equations, ACM TOMS, 2015.



Outline



Basic Adjoint AD

First Order Implementation Case Study

Beyond Basic Adjoint AD

Checkpointing Preaccumulation Higher-Level Adjoint Elementals Debugging Adjoints Further Aspects

Conclusion

Recent Applications / Publications

Miscellaneous

Inside 1st-Order Algorithmic Adjoints Beyond 1st-Order Adjoints Further Details



Adjoint interpretation (of DAG) yields

 $\mathbf{p}_{(1)} = \mathbf{p}_{(1)}(\mathbf{p}, \mathbf{x}_{(1)})$,

e.g, with $v_{i(1)} = 0$ for i = 2, 3, 4

 $\begin{aligned} & v_{4(1)} += 42 \cdot x_{1(1)} \\ & v_{4(1)} += p_0 \cdot x_{0(1)} \\ & p_{0(1)} += v_4 \cdot x_{0(1)} \\ & \cdots \end{aligned}$

Edges (local partial derivatives; PMR = |E|) are accessed sequentially.

Vertices (adjoints; PMR = |V|) are accessed "randomly."





Problem Description and Summary of Results [6]

- infeasible memory requirement of naive adjoints
- automatically, stream DAG to slower³ storage instead of nontrivial checkpointing
- "page faults" due to random access to adjoints
- Idea: separate sequentially from randomly accessed data; stream sequentially accessed data to slower storage; keep randomly accessed data in fast storage \Rightarrow minimize adjoint memory requirement (AMR)

$\mathsf{AMR} \stackrel{?}{\ll} |V|$

- AMR = # program variables (\rightarrow source transformation)
- ► AMR = $\max_{(i,j)\in E}(j-i)$ ⇒ DIRECTED BANDWIDTH
- decomposition by perpetuation
- case study: Discrete Adjoint OpenFOAM

³speed of access

Reduction of Adjoint Memory Requirement E.g. Variant of $Bat DAG \Rightarrow AMR = 12$





Bat [N. 2004]



... $\% \max_{(i,j) \in E} (j-i) \Rightarrow \mathsf{AMR} = 6$





... DIRECTED BANDWIDTH $\Rightarrow AMR = 5$





.. Perpetuation \Rightarrow AMR = 3



Sensitivity of pressure loss from inlet (left) to outlet (right) over 122,250 cells.



 $^a\mathrm{Run}$ time penalty <10%

Beyond 1st-Order Algorithmic Adjoints E.g. Hessian by $\mathbf{p}_{(1)}^{(2)} += x_{(1)} \cdot \nabla^2 f(\mathbf{p}) \cdot \mathbf{p}^{(2)}$



```
1
   template<typename AFLT_T, typename PFLT_T>
   void ddxdpdp_t(AFLT_T& x_v, std::vector<AFLT_T>& p_v,
2
            const std::vector<std::vector<PFLT T>> &dW.
3
            std::vector<std::vector<AFLT_T>> &ddxdpp) {
4
      typedef typename dco::gt1s<AFLT_T>::type DCO_T;
5
      int n=p v.size();
6
      std::vector<DCO_T> p(n), dxdp(n); dco::value(p)=p_v;
7
      for (int i=0;i<n;i++) {</pre>
8
        DCO T x=x v:
q
        dco::derivative(p[i])=1;
10
        dxdp_a(x,p,dW,dxdp); // first-order adjoint driver
        ddxdpp[i]=dco::derivative(dxdp);
12
        dco::derivative(p[i])=0;
13
     7
14
15
   }
```

dco/c++ allows for seamless instantiation of arbitrary combinations of tangent and (non-black-box) adjoint modes.

Beyond 1st-Order Algorithmic Adjoints E.g. Hessian by $\mathbf{p}_{(2)}^{(1)} += x_{(2)}^{(1)} \cdot \nabla^2 f(\mathbf{p}) \cdot \mathbf{p}^{(1)}$



```
template<typename AFLT_T, typename PFLT_T>
1
   void ddxdpdp_a(
2
            AFLT_T& x_v, std::vector<AFLT_T>& p_v,
3
            const std::vector<std::vector<PFLT_T>> &dW,
            std::vector<std::vector<AFLT_T>> &ddxdpdp
5
            ) {
6
      typedef typename dco::gals<AFLT T> DCO M:
7
      typedef typename DCO_M::type DCO_T;
8
      typedef typename DCO_M::tape_t DCO_TT;
q
      int n=p_v.size(); std::vector<DCO_T> p(n); dco::value(p)=p_v;
10
      for (int i=0;i<n;i++) {</pre>
11
        DCO T x=x v:
12
        std::vector<DCO_T> dxdp(n,0);
13
        DCO_M::global_tape=DCO_TT::create();
14
        DCO_M::global_tape->register_variable(p);
15
        dxdp_t(x,p,dW,dxdp); // first-order tangent driver
16
        dco::derivative(dxdp[i])=1;
17
        DCO_M::global_tape->interpret_adjoint();
18
        ddxdpdp[i]=dco::derivative(p);
19
        DCO_TT::remove(DCO_M::global_tape);
20
      }
21
22
   }
```

Beyond 1st-Order Algorithmic Adjoints E.g. Hessian by $\mathbf{p}_{(2)} += x_{(1)} \cdot \nabla^2 f(\mathbf{p}) \cdot \mathbf{p}_{(1,2)}$



```
1
   template<typename AFLT T, typename PFLT T>
   void ddxdpdp_a(
2
            AFLT_T& x_v, std::vector<AFLT_T>& p_v,
3
            const std::vector<std::vector<PFLT T>> &dW.
4
            std::vector<std::vector<AFLT_T>> &ddxdpdp
5
            ) {
6
      typedef typename dco::gals<AFLT_T> DCO_M;
7
      typedef typename DCO_M::type DCO_T;
8
      typedef typename DCO_M::tape_t DCO_TT;
q
      int n=p_v.size(); std::vector<DC0_T> p(n); dco::value(p)=p_v;
10
      for (int i=0:i<n:i++) {</pre>
11
        DCO_T x=x_v;
12
        std::vector<DCO_T> dxdp(n,0);
13
        DCO_M::global_tape=DCO_TT::create();
14
        DCO_M::global_tape->register_variable(p);
15
        dxdp_a(x,p,dW,dxdp); // first-order adjoint driver
16
        dco::derivative(dxdp[i])=1;
17
        DCO_M::global_tape->interpret_adjoint();
18
        ddxdpdp[i]=dco::derivative(p);
19
        DCO_TT::remove(DCO_M::global_tape);
20
     }
21
22
```



The tangent

$$A^{(1)} \cdot \mathbf{x} + A \cdot \mathbf{x}^{(1)} = \mathbf{b}^{(1)}$$

holds at the primal solution ${\bf x}$ (see tangent matrix-vector product).

The corresponding adjoint

$$A_{(1)} = \mathbf{b}_{(1)} \cdot \mathbf{x}^T$$
$$\mathbf{x}_{(1)} = A^T \cdot \mathbf{b}_{(1)}$$

follows from

$$Y^{(1)} = \sum_{i} A_i \cdot X_i^{(1)} \cdot B_i \quad \Leftrightarrow \quad X_{i(1)} = A_i^T \cdot Y_{(1)} \cdot B_i^T .$$

[3]





$$\frac{dF}{d\mathbf{p}}(\mathbf{x},\mathbf{p})\cdot\mathbf{p}^{(1)}=0$$

implies at the primal solution $\mathbf{x}=\mathbf{x}(\mathbf{p})$

$$\mathbf{x}^{(1)} = -\frac{dF}{d\mathbf{x}}^{-1} \cdot \frac{\partial F}{\partial \mathbf{p}} \cdot \mathbf{p}^{(1)}$$

The adjoint

$$\mathbf{p}_{(1)} = -\frac{\partial F}{\partial \mathbf{p}}^T \cdot \frac{dF}{d\mathbf{x}}^{-T} \cdot \mathbf{x}_{(1)}$$

follows from

$$Y^{(1)} = \sum_{i} A_i \cdot X_i^{(1)} \cdot B_i \quad \Leftrightarrow \quad X_{i(1)} = A_i^T \cdot Y_{(1)} \cdot B_i^T .$$

[3]

Naumann, Paris, 24.01.2019

 $\min_{\mathbf{x}} f(\mathbf{x}(\mathbf{p}), \mathbf{p})$ s.t. $g(\mathbf{x}(\mathbf{p}), \mathbf{p}) = 0$ Further Details I



KKT

$$0 = \frac{d\mathcal{L}}{d(\mathbf{x},\lambda)}(\mathbf{x},\mathbf{p},\lambda) = \frac{d\left(f(\mathbf{x},\mathbf{p}) + \lambda^T \cdot g(\mathbf{x},\mathbf{p})\right)}{d(\mathbf{x},\lambda)}$$

implies

$$\frac{d\left(f(\mathbf{x}, \mathbf{p}) + \lambda^T \cdot g(\mathbf{x}, \mathbf{p})\right)}{d\mathbf{x}} = 0$$
$$g(\mathbf{x}, \mathbf{p}) = 0$$

at a solution $(\mathbf{x},\lambda)=(\mathbf{x}(\mathbf{p}),\lambda(\mathbf{p})).$





Hence, at the solution

$$\begin{aligned} \frac{d\mathcal{L}}{d\mathbf{p}}(\mathbf{x},\mathbf{p},\lambda) &= \frac{df}{d\mathbf{p}}(\mathbf{x},\mathbf{p}) \\ &= \frac{df}{d\mathbf{x}} \cdot \frac{d\mathbf{x}}{d\mathbf{p}} + \frac{\partial f}{\partial \mathbf{p}} + \lambda^T \cdot \left(\frac{dg}{d\mathbf{x}} \cdot \frac{d\mathbf{x}}{d\mathbf{p}} + \frac{\partial g}{\partial \mathbf{p}}\right) \\ &= \underbrace{\left(\frac{df}{d\mathbf{x}} + \lambda^T \cdot \frac{dg}{d\mathbf{x}}\right)}_{=0} \cdot \frac{d\mathbf{x}}{d\mathbf{p}} + \frac{\partial f}{\partial \mathbf{p}} + \lambda^T \cdot \frac{\partial g}{\partial \mathbf{p}} \end{aligned}$$

implying

$$\frac{df}{d\mathbf{p}}(\mathbf{x}) = \lambda^T \cdot \frac{\partial g}{\partial \mathbf{p}}$$

as in this case

$$\frac{\partial f}{\partial \mathbf{p}} = 0 \; .$$