

AD Master Class: Advanced Adjoint Techniques

Computing Hessians

Viktor Mosenkis
viktor.mosenkis@nag.co.uk



Experts in numerical algorithms
and HPC services

25 November 2020

AD Masterclass Schedule and Remarks

■ AD Masterclass Schedule

- 1 October 2020 | Checkpointing and external functions 1
- 15 October 2020 | Checkpointing and external functions 2
- 29 October 2020 | Guest lecture by Prof Uwe Naumann on Advanced AD topics in Machine Learning
- 12 November 2020 | Monte Carlo
- 19 November 2020 | Guest lecture by Prof Uwe Naumann on Adjoint Code Design Patterns applied to Monte Carlo
- 25 November 2020 | Computing Hessians

■ Remarks

- Please submit your questions via the questions panel at any time during this session, these will be addressed at the end.
- A recording of this session, along with the slides will be shared with you in a day or two.

Dialogue

We want this webinar series to be interactive (even though it's hard to do)

- We want your feedback, we want to adapt material to your feedback
- Please feel free to contact us via email to ask questions at any time
- We'd love to reach out offline, discuss what's working, what to spend more time on
- For some orgs, may make sense for us to do a few bespoke sessions

- This is an advanced course
- We assume that you are familiar with the material from the first Masterclass series
- You will get access to the materials from the first Masterclass series via email in a day or two
- Also it is not a pre-requisite we recommend to review the material from the previous series
- We will try to give references to the previous Masterclass series whenever possible

Outcomes

- Discuss four different second-order models
 - their advantages and disadvantages
 - usage of the models to compute the full Hessian
- How to reuse external function implemented for the first order adjoint code in the second order model

Second-Order Models

To compute second-order derivatives we can simply apply tangent or adjoint model to tangent/adjoint model, yielding the following four second-order models

- tangent over tangent (forward over forward)
- tangent over adjoint (forward over reverse)
- adjoint over tangent (reverse over forward)
- adjoint over adjoint (reverse over reverse)

Tangent over Tangent Model (Second-Order Tangent)

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

A second-order tangent code

$$\tilde{F} : \mathbb{R}^n \times \mathbb{R}^n \times \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}^m \times \mathbb{R}^m \times \mathbb{R}^m \times \mathbb{R}^m,$$

$$(\mathbf{y}, \tilde{\mathbf{y}}, \dot{\mathbf{y}}, \tilde{\dot{\mathbf{y}}})^T = \tilde{F}(\mathbf{x}, \tilde{\mathbf{x}}, \dot{\mathbf{x}}, \tilde{\dot{\mathbf{x}}})$$

computes a mixture of first and second derivative information alongside with the function value as follows:

$$\mathbf{y} = F(\mathbf{x})$$

$$\tilde{\mathbf{y}} = F'(\mathbf{x}) \cdot \tilde{\mathbf{x}}$$

$$\dot{\mathbf{y}} = F'(\mathbf{x}) \cdot \dot{\mathbf{x}}$$

$$\tilde{\dot{\mathbf{y}}} = \dot{\mathbf{x}}^T \cdot F''(\mathbf{x}) \cdot \tilde{\mathbf{x}} + F'(\mathbf{x}) \cdot \tilde{\mathbf{x}}$$

Tangent over Tangent: Accumulation of Hessian $m = 1$

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

As $m = 1$ $F''(\mathbf{x})$ is an n by n matrix.

In the tangent over tangent model the Hessian can be extracted from the computation of

$$\tilde{\mathbf{y}} = \underbrace{\dot{\mathbf{x}}^T}_{\in \mathbb{R}^n} \cdot F''(\mathbf{x}) \cdot \underbrace{\tilde{\mathbf{x}}}_{\in \mathbb{R}^n} + F'(\mathbf{x}) \cdot \tilde{\mathbf{x}}$$

Setting $\dot{\mathbf{x}} = \mathbf{e}_i$ and $\tilde{\mathbf{x}} = \mathbf{e}_j$, while keeping $\tilde{\mathbf{x}} = \mathbf{0}$ yields

$\tilde{\mathbf{y}} = F''(\mathbf{x})_{i,j}$ (the i, j -th entry of the Hessian).

Hence computation of Hessian with tangent over tangent is $O(n^2)$ as we need to range over the Cartesian basis vectors of \mathbb{R}^n in both $\dot{\mathbf{x}}$ and $\tilde{\mathbf{x}}$.

Hessian matrix is symmetric so we need to compute only upper or lower triangular part of the Hessian

Tangent over Tangent: Implementation

In dco/c++ tangent over tangent is implemented via

```
1 using DCO_BASE_M = dco::gt1s<double>;
2 using DCO_BASE_T = DCO_BASE_M::type;
3 using DCO_M = dco::gt1s<DCO_BASE_T>;
```

- $x \rightarrow \text{dco::value}(\text{dco::value}(x))$
- $\dot{x} \rightarrow \text{dco::value}(\text{dco::derivative}(x))$
- $\tilde{x} \rightarrow \text{dco::derivative}(\text{dco::value}(x))$
- $\tilde{\dot{x}} \rightarrow \text{dco::derivative}(\text{dco::derivative}(x))$

Tangent over Tangent: Driver

```
1  for (size_t i = 0; i < n; i++) {
2      for (size_t j = 0; j < n; j++) {
3          dco::value(dco::derivative(x[i])) = 1.0; //\dot{x}
4          dco::derivative(dco::value(x[j])) = 1.0; //\tilde{x}
5
6          foo(n, x, y);
7          //H[i][j] = \tilde{\dot{y}}
8          Hess[i][j] = dco::derivative(dco::derivative(y));
9
10         dco::value(dco::derivative(x[i])) = 0.0; //\dot{x}
11         dco::derivative(dco::value(x[j])) = 0.0; //\tilde{x}
12     }
13 }
```

Tangent over Adjoint Model (Second-Order Adjoint)

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

A second-order adjoint code

$$\dot{\bar{F}} : \mathbb{R}^n \times \mathbb{R}^n \times \mathbb{R}^m \times \mathbb{R}^m \rightarrow \mathbb{R}^m \times \mathbb{R}^m \times \mathbb{R}^n \times \mathbb{R}^n,$$

$$(\mathbf{y}, \dot{\mathbf{y}}, \bar{\mathbf{x}}, \dot{\bar{\mathbf{x}}})^T = \dot{\bar{F}}(\mathbf{x}, \dot{\mathbf{x}}, \bar{\mathbf{y}}, \dot{\bar{\mathbf{y}}})$$

computes a mixture of first and second derivative information alongside with the function value as follows:

$$\mathbf{y} = F(\mathbf{x})$$

$$\dot{\mathbf{y}} = F'(\mathbf{x}) \cdot \dot{\mathbf{x}}$$

$$\bar{\mathbf{x}} = F'(\mathbf{x})^T \cdot \bar{\mathbf{y}}$$

$$\dot{\bar{\mathbf{x}}} = \bar{\mathbf{y}}^T \cdot F''(\mathbf{x}) \cdot \dot{\mathbf{x}} + F'(\mathbf{x})^T \cdot \dot{\bar{\mathbf{y}}}$$

Tangent over Adjoint: Accumulation of Hessian $m = 1$

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

As $m = 1$ $F''(\mathbf{x})$ is an n by n matrix.

In the tangent over adjoint model the Hessian can be extracted from the computation of

$$\dot{\mathbf{x}} = \underbrace{\bar{\mathbf{y}}^T \cdot F''(\mathbf{x}) \cdot \dot{\mathbf{x}}}_{\in \mathbb{R}} + F'(\mathbf{x})^T \cdot \dot{\mathbf{y}}$$

Setting $\bar{\mathbf{y}} = \mathbf{1}$ and $\dot{\mathbf{x}} = \mathbf{e}_i$, while keeping $\dot{\mathbf{y}} = \mathbf{0}$ yields $\dot{\mathbf{x}} = F''(\mathbf{x})_i$ (the i -th row (or column) of the Hessian due to symmetry).

Hence computation of Hessian with tangent over adjoint is $O(n)$ as we need to range over the Cartesian basis vectors of \mathbb{R}^n in $\dot{\mathbf{x}}$.

Tangent over Adjoint: Implementation

In dco/c++ tangent over adjoint is implemented via

```
1 using DCO_BASE_M = dco::gt1s<double>;
2 using DCO_BASE_T = DCO_BASE_M::type;
3 using DCO_M = dco::ga1s<DCO_BASE_T>;
```

- $x \rightarrow \text{dco::value}(\text{dco::value}(x))$
- $\bar{x} \rightarrow \text{dco::value}(\text{dco::derivative}(x))$
- $\dot{x} \rightarrow \text{dco::derivative}(\text{dco::value}(x))$
- $\dot{\bar{x}} \rightarrow \text{dco::derivative}(\text{dco::derivative}(x))$

Tangent over Adjoint: Driver

```
1  for (size_t k = 0; k < n; k++) {  
2      DCO_M::global_tape = DCO_TAPE_T::create(o);  
3  
4      DCO_M::global_tape->register_variable(x);  
5      dco::derivative(dco::value(x[k])) = 1.0; //\dot{x}_i = 1  
6  
7      foo(n, x, y); // Record the tape  
8  
9      dco::value(dco::derivative(y)) = 1.0; //\bar{Y}  
10  
11     DCO_M::global_tape->interpret_adjoint();  
12     for (int i = 0; i < n; i++)  
13         Hess[k][i] = dco::derivative(dco::derivative(x[i]));  
14  
15     dco::derivative(dco::value(x[k])) = 0.0; //\dot{x}_i = 0  
16     DCO_TAPE_T::remove(DCO_M::global_tape);  
17 }
```

Adjoint over Tangent Model

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

A second-order adjoint over tangent model

$$\dot{\bar{F}} : \mathbb{R}^n \times \mathbb{R}^n \times \mathbb{R}^m \times \mathbb{R}^m \rightarrow \mathbb{R}^m \times \mathbb{R}^m \times \mathbb{R}^n \times \mathbb{R}^n,$$

$$(\mathbf{y}, \dot{\mathbf{y}}, \bar{\mathbf{x}}, \dot{\bar{\mathbf{x}}})^T = \dot{\bar{F}}(\mathbf{x}, \dot{\mathbf{x}}, \bar{\mathbf{y}}, \dot{\bar{\mathbf{y}}})$$

computes a mixture of first and second derivative information alongside with the function value as follows:

$$\mathbf{y} = F(\mathbf{x})$$

$$\dot{\mathbf{y}} = F'(\mathbf{x}) \cdot \dot{\mathbf{x}}$$

$$\bar{\mathbf{x}} = F'(\mathbf{x})^T \cdot \bar{\mathbf{y}} + \dot{\bar{\mathbf{y}}}^T \cdot F''(\mathbf{x}) \cdot \dot{\mathbf{x}}$$

$$\dot{\bar{\mathbf{x}}} = F'(\mathbf{x})^T \cdot \bar{\mathbf{y}}$$

Adjoint over Tangent: Accumulation of Hessian $m = 1$

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

As $m = 1$ $F''(\mathbf{x})$ is an n by n matrix.

In the adjoint over tangent model the Hessian can be extracted from the computation of

$$\bar{\mathbf{x}} = F'(\mathbf{x})^T \cdot \bar{\mathbf{y}} + \underbrace{\bar{\mathbf{y}}^T \cdot F''(\mathbf{x}) \cdot \dot{\mathbf{x}}}_{\in \mathbb{R}^n}$$

Setting $\bar{\mathbf{y}} = \mathbf{1}$ and $\dot{\mathbf{x}} = \mathbf{e}_i$, while keeping $\bar{\mathbf{y}} = \mathbf{0}$ yields $\bar{\mathbf{x}} = F''(\mathbf{x})_i$ (the i -th row (or column) of the Hessian due to symmetry).

Hence computation of Hessian with tangent over adjoint is $O(n)$ as we need to range over the Cartesian basis vectors of \mathbb{R}^n in $\dot{\mathbf{x}}$.

Adjoint over Tangent: Implementation

In dco/c++ adjoint over tangent is implemented via

```
1 using DCO_BASE_M = dco::ga1s<double>;
2 using DCO_BASE_T = DCO_BASE_M::type;
3 using DCO_M = dco::gt1s<DCO_BASE_T>;
```

- $x \rightarrow \text{dco::value}(\text{dco::value}(x))$
- $\dot{x} \rightarrow \text{dco::value}(\text{dco::derivative}(x))$
- $\bar{x} \rightarrow \text{dco::derivative}(\text{dco::value}(x))$
- $\dot{\bar{x}} \rightarrow \text{dco::derivative}(\text{dco::derivative}(x))$

Adjoint over Tangent: Driver

```
1  using B_M = dco::ga1s<double>;
2  for (size_t k = 0; k < n; k++) {
3      B_M::global_tape = DCO_TAPE_T::create(o);
4      //\bar{x}
5      B_M::global_tape->register_variable(dco::value(x));
6      dco::value(dco::derivative(x[k])) = 1.0; //\dot{x}_i
7
8      foo(n, x, y);
9      dco::derivative(dco::derivative(y)) = 1.0; //\bar{\dot{y}}
10
11     B_M::global_tape->interpret_adjoint();
12
13     for (int i = 0; i < n; i++) //\bar{x}
14         Hess[k][i] = dco::derivative(dco::value(x[i]));
15     dco::value(dco::derivative(x[k])) = 0.0; //\dot{x}_i = 0
16     DCO_TAPE_T::remove(B_M::global_tape);
17 }
```

Adjoint over Adjoint Model

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

A second-order adjoint over adjoint model

$$\tilde{\tilde{F}} : \mathbb{R}^n \times \mathbb{R}^m \times \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}^m \times \mathbb{R}^n \times \mathbb{R}^m \times \mathbb{R}^n,$$

$$(\mathbf{y}, \bar{\mathbf{x}}, \tilde{\bar{\mathbf{y}}}, \tilde{\bar{\mathbf{x}}})^T = \tilde{\tilde{F}}(\mathbf{x}, \bar{\mathbf{y}}, \tilde{\bar{\mathbf{x}}}, \tilde{\bar{\mathbf{y}}})$$

computes a mixture of first and second derivative information alongside with the function value as follows:

$$\mathbf{y} = F(\mathbf{x})$$

$$\bar{\mathbf{x}} = F'(\mathbf{x})^T \cdot \bar{\mathbf{y}}$$

$$\tilde{\bar{\mathbf{y}}} = F'(\mathbf{x}) \cdot \tilde{\bar{\mathbf{x}}}$$

$$\tilde{\bar{\mathbf{x}}} = F'(\mathbf{x})^T \cdot \tilde{\bar{\mathbf{y}}} + \bar{\mathbf{y}}^T \cdot F''(\mathbf{x}) \cdot \tilde{\bar{\mathbf{x}}}$$

Adjoint over Adjoint: Accumulation of Hessian $m = 1$

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

As $m = 1$ $F''(\mathbf{x})$ is an n by n matrix.

In the adjoint over adjoint model the Hessian can be extracted from the computation of

$$\tilde{\mathbf{x}} = F'(\mathbf{x})^T \cdot \tilde{\mathbf{y}} + \underbrace{\tilde{\mathbf{y}}^T}_{\in \mathbb{R}} \cdot F''(\mathbf{x}) \cdot \underbrace{\tilde{\mathbf{x}}}_{\in \mathbb{R}^n}$$

Setting $\tilde{\mathbf{y}} = 1$ and $\tilde{\mathbf{x}} = e_i$, while keeping $\tilde{\mathbf{y}} = 0$ yields $\tilde{\mathbf{x}} = F''(\mathbf{x})_i$ (the i -th row (or column) of the Hessian due to symmetry).

Hence computation of Hessian with tangent over adjoint is $O(n)$ as we need to range over the Cartesian basis vectors of \mathbb{R}^n in $\tilde{\mathbf{x}}$.

No advantage using the adjoint model for the second time

Why there is no advantage in using adjoint model twice?

The original function is

$$F : \mathbb{R}^n \rightarrow \mathbb{R}, \quad \mathbf{y} = F(\mathbf{x}) = F(x_1, \dots, x_n)$$

The corresponding derivative function (that can be computed by one evaluation of the first order adjoint model) is

$$G : \mathbb{R}^n \rightarrow \mathbb{R}^n, \quad (y_1, \dots, y_n)^T = \left(\frac{\partial F(\mathbf{x})}{\partial x_1}, \dots, \frac{\partial F(\mathbf{x})}{\partial x_n} \right)$$

The second application of the adjoint model we differentiate G . The input and output dimension of G are equal, therefore it doesn't matter what models is used for the second time.

Adjoint over Adjoint: Implementation

In dco/c++ adjoint over tangent is implemented via

```
1 using DCO_BASE_M = dco::ga1s<double>;
2 using DCO_BASE_T = DCO_BASE_M::type;
3 using DCO_M = dco::ga1s<DCO_BASE_T>;
```

- $x \rightarrow \text{dco::value}(\text{dco::value}(x))$
- $\bar{x} \rightarrow \text{dco::value}(\text{dco::derivative}(x))$
- $\tilde{x} \rightarrow \text{dco::derivative}(\text{dco::value}(x))$
- $\tilde{\bar{x}} \rightarrow \text{dco::derivative}(\text{dco::derivative}(x))$

Adjoint over Adjoint: Driver

```
1  using B_M = dco::ga1s<double>;
2
3  B_M::global_tape = DCO_BASE_TAPE_T::create(o);
4  DCO_M::global_tape = DCO_TAPE_T::create(o);
5
6  for (size_t i = 0; i < n; ++i) {
7      //need \tilde{x}
8      B_M::global_tape->register_variable(dco::value(x[i]));
9      //need to differentiate \bar{x} to compute the adjoints
10     DCO_M::global_tape->register_variable(x[i]);
11 }
12
13 //record both tapes
14 foo(n, x, y);
```

Adjoint over Adjoint: Driver

```
1  for (size_t k = 0; k < n; k++) {  
2      dco::value(dco::derivative(y)) = 1.0; //\bar{y} = 1.0  
3      DCO_M::global_tape->interpret_adjoint();  
4  
5      //\tilde{\bar{x}} = e_k  
6      dco::derivative(dco::derivative(x[k])) = 1.0;  
7  
8      B_M::global_tape->interpret_adjoint();  
9  
10     //Extract the derivatives  
11     for (int i = 0; i < n; i++) //\tilde{x}  
12         Hess[k][i] = dco::derivative(dco::value(x[i]));  
13  
14     DCO_M::global_tape->zero_adjoint();  
15     B_M::global_tape->zero_adjoint();  
16 }
```

Summary Second-Order Models ($m = 1$)

- Second-Order tangent model has
 - smallest memory requirements
 - $O(n^2)$ complexity for accumulating the Hessian
- All three adjoint models
 - are mathematically equivalent ($O(n)$ for the Hessian)
 - implementation differ
- No advantage using the adjoint model for the second time

Summary Second-Order Adjoint Models ($m = 1$)

■ tangent over adjoint

- the preferred second-order adjoint model
- smallest tape size from all adjoint models
- external functions implementation from the first order model can be reused

■ adjoint over tangent

- higher tape size compared to tangent over adjoint

■ adjoint over adjoint

- highest memory requirements for the tape
- complicated drivers due to two tapes
- no re-recording of the tapes needed for computing Hessian projections
(can be faster compared second order models involving tangent computation in some cases)

Second-Order Models ($m > 1$)

For $m > 1$ the Hessian $F''(\mathbf{x})$ is a $m \times n \times n$ tensor.

- tangent over tangent $\tilde{\mathbf{y}} = \underbrace{\dot{\mathbf{x}}^T}_{\in \mathbb{R}^n} \cdot F''(\mathbf{x}) \cdot \underbrace{\tilde{\mathbf{x}}}_{\in \mathbb{R}^n} + F'(\mathbf{x}) \cdot \tilde{\mathbf{x}}$. For each $\dot{\mathbf{x}} = e_i$, $\tilde{\mathbf{x}} = e_j$ we compute i, j -th entry in all $n \times n$ sub-matrices of F'' tensor. Hence $O(n^2)$ complexity.
- tangent over adjoint $\dot{\mathbf{x}} = \underbrace{\bar{\mathbf{y}}^T}_{\in \mathbb{R}^m} \cdot F''(\mathbf{x}) \cdot \underbrace{\dot{\mathbf{x}}}_{\in \mathbb{R}^n} + F'(\mathbf{x})^T \cdot \dot{\mathbf{y}}$ For each $\dot{\mathbf{x}} = e_j$, $\bar{\mathbf{y}} = e_i$ we compute i, j -th entry in all $m \times n$ sub-matrices of F'' tensor. Hence $O(m \cdot n)$ complexity.

Second-Order Models ($m > 1$)

For $m > 1$ the Hessian $F''(\mathbf{x})$ is a $m \times n \times n$ tensor.

- adjoint over tangent $\bar{\mathbf{x}} = F'(\mathbf{x})^T \cdot \bar{\mathbf{y}} + \underbrace{\bar{\mathbf{y}}^T \cdot F''(\mathbf{x}) \cdot \dot{\mathbf{x}}}_{\in \mathbb{R}^m \quad \in \mathbb{R}^n}$ For each $\dot{\mathbf{x}} = e_j$, $\bar{\mathbf{y}} = e_i$ we compute i, j -th entry in all $m \times n$ sub-matrices of F'' tensor. Hence $O(m \cdot n)$ complexity.
- adjoint over adjoint $\tilde{\mathbf{x}} = F'(\mathbf{x})^T \cdot \tilde{\mathbf{y}} + \underbrace{\tilde{\mathbf{y}}^T \cdot F''(\mathbf{x}) \cdot \tilde{\bar{\mathbf{x}}}}_{\in \mathbb{R}^m \quad \in \mathbb{R}^n}$ For each $\tilde{\bar{\mathbf{x}}} = e_j$, $\bar{\mathbf{y}} = e_i$ we compute i, j -th entry in all $m \times n$ sub-matrices of F'' tensor. Hence $O(m \cdot n)$ complexity.

Reusing first order adjoint code for external functions

Summary

In this Masterclass we

- learned the four different second-order models
- learned how to use the second-order models to compute the Hessian
- discussed the complexity of computing the Hessian with different models
- learned how to reuse external function implementation from the first-order adjoint model in the second-order adjoint model

You will see a survey on your screen after exiting
from this session.

We would appreciate your feedback.

We are now moving on the Q&A Session