# Introduction to Algorithmic Differentiation

## AD by Hand (Tangent Code)

Uwe Naumann

Informatik 12:
Software and Tools for Computational Engineering (STCE)

RWTH Aachen

# Contents

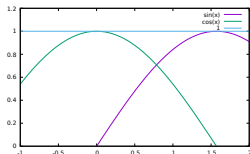Let $f : \mathbf{R}^2 \to \mathbf{R}$ be defined as

$$f(x, p) = \begin{cases} f_1(x) & x < p \\ f_2(x) & x \geq p \, . \end{cases}$$



with differentiable univariate scalar $f_1$ and $f_2$.

Depending on the choice of $f_1$ and $f_2$ the function $f$ can be nondifferentiable or even discontinuous at $x = p$.

Examples:

▶ $f_1 = \cos$, $f_2 = \sin \Rightarrow$ discontinuous at $x = p = 1$
▶ $f_1 = \cos$, $f_2 = \sin \Rightarrow$ nondifferentiable at $x = p = \frac{\pi}{4}$
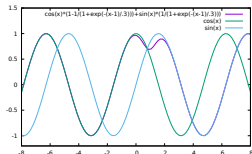▶ $f_1 = 1$, $f_2 = \cos \Rightarrow$ differentiable at $x = p = 0$

Sigmoidal smoothing replaces $f$ with $\tilde{f} : \mathbf{R}^3 \to \mathbf{R}$ defined as

$$\tilde{f}(x, p, w) = (1 - \sigma(x, p, w)) \cdot f_1(x) + \sigma(x, p, w) \cdot f_2(x) ,$$
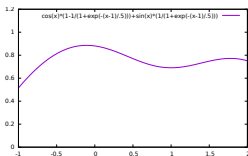
where

$$\sigma(x, p, w) = \frac{1}{1 + e^{-\frac{x - p}{w}}} .$$



$w = 0.3$

Example: $f_1 = \cos$, $f_2 = \sin$ at $x = p = 1$



$w = 0.5$



$w = 0.1$



$w = 0.05$

# Sigmoidal Smoothing
## Implementation

nondifferentiable $\Rightarrow$ differentiable

$\Rightarrow$

```
1   #pragma once
2
3   #include <cmath>
4
5   template<typename T>
6   void g(T &x, const T &p)
       {
7     if (x<p)
8       x=sin(x);
9     else
10      x=cos(x);
11  }
```

```
1   #pragma once
2
3   #include <cmath>
4
5   template<typename T>
6   void f(T &x, const T &p, const T &w) {
7     T f1=sin(x);
8     T f2=cos(x);
9     x=1./(1.+exp(-(x-p)/w));
10    x=f1*(1-x)+f2*x;
11  }
```

Consider a nonlinear equation $y = f(x) = 0$ at some (starting) point $x$.

Building on the assumption that $f(x + \Delta x) \approx f(x) + f'(x) \cdot \Delta x$ the root finding problem for $f$ can be replaced locally by the root finding problem for the linearization

$$\bar{f}(\Delta x) = f(x) + f'(x) \cdot \Delta x \ .$$

The right-hand side is a straight line intersecting the $y$-axis in $(\Delta x = 0, \bar{f}(\Delta x) = f(x))$.

Solution of

$$\bar{f}(\Delta x) = f(x) + f'(x) \cdot \Delta x = 0$$

for $\Delta x$ yields

$$\Delta x = -\frac{f(x)}{f'(x)}$$

implying $f(x + \Delta x) \approx 0$.

If the new iterate is not close enough to the root of the nonlinear function, i.e, $|f(x + \Delta x)| > \epsilon$ for some measure of accuracy of the numerical approximation $\epsilon > 0$, then it becomes the starting point for the next iteration yielding the recurrence

$$x = x - \frac{f(x)}{f'(x)}$$

Convergence of this method is not guaranteed in general. Damping of the magnitude of the next step may help.

$$x = x - \alpha \cdot \frac{f(x)}{f'(x)} \quad \text{for } 0 < \alpha \leq 1 .$$

The damping parameter $\alpha$ is often determined by line search (e.g, recursive bisection yielding $\alpha = 1, \ 0.5, \ 0.25, \ \ldots$) such that decrease in absolute function value is ensured.

# Newton's Method
## Implementation (E.g., $x^2 - p = 0$)

The following iteration terminates if the residual is close enough to zero or if a given number (maxit) of iterations was performed.

```
1  template<typename T, typename PT>
2  void newton(T &x, const T &p, const PT
       &eps, const unsigned int maxit) {
3    unsigned int it=0;
4    T f=pow(x,2)−p;
5    do {
6      T dfdx=2∗x;
7      x−=f/dfdx;
8      f=pow(x,2)−p;
9      if (++it==maxit) break;
10   } while(fabs(f)>eps);
11 }
```

Alternatively, …

```
1  template<typename T, typename PT>
2  T f(T &x, const PT &p) {
3    return pow(x,2)−p;
4  }
5
6  template<typename T, typename PT>
7  T dfdx(T &x, const PT &) { return 2∗x; }
8
9  template<typename T, typename PT>
10 void newton(T &x, const T &p, const PT
       &eps, const unsigned int maxit) {
11   unsigned int it=0;
12   T y=f(x,p);
13   do {
14     x−=y/dfdx(x,p);
15     y=f(x,p);
16     if (++it==maxit) break;
17   } while(fabs(y)>eps);
18 }
```

We consider differentiable numerical programs

$$\begin{pmatrix} y \\ \tilde{y} \end{pmatrix} = F(x, \tilde{x}) : \boldsymbol{R}^n \times \boldsymbol{R}^{\tilde{n}} \to \boldsymbol{R}^m \times \boldsymbol{R}^{\tilde{m}}$$

mapping active ($x \in \boldsymbol{R}^n$) and passive ($\tilde{x} \in \boldsymbol{R}^{\tilde{n}}$) inputs onto active ($y \in \boldsymbol{R}^m$) and passive ($\tilde{y} \in \boldsymbol{R}^{\tilde{m}}$) outputs. The active output $y$ is assumed to be differentiable with respect to $x$ with

$$F' \equiv F'(x, \tilde{x}) = \frac{dy}{dx} \ .$$

The corresponding (first-order) tangent program computes

$$\begin{pmatrix} y \\ \tilde{y} \\ y^{(1)} \end{pmatrix} = F^{(1)}(x, \tilde{x}, x^{(1)}) \equiv \begin{pmatrix} F(x, \tilde{x}) \\ F' \cdot x^{(1)} \end{pmatrix} \ .$$

# Tangent Programs
## Example: Sigmoidal Smoothing

▶ all arguments active for $\frac{dx}{d(x\ p\ w)^T}$ or $\frac{dx}{d(p\ w)^T}$

```
1  template<typename T>
2  void f_t(T &x, T &x_t, const T &p, const T &p_t, const T &w, const T &w_t) {
3    ...
4  }
```

▶ p passive for $\frac{dx}{d(x\ w)^T}$ or $\frac{dx}{dw}$

```
1  template<typename T>
2  void f_t(T &x, T &x_t, const T &p, const T &w, const T &w_t) {
3    ...
4  }
```

▶ w passive for $\frac{dx}{d(x\ p)^T}$ or $\frac{dx}{dp}$

```
1  template<typename T>
2  void f_t(T &x, T &x_t, const T &p, const T &p_t, const T &w) {
3    ...
4  }
```

For given values of the inputs $x = (x_i)_{i=0}^{n-1}$ and $\tilde{x}$ the active section

$$y = (y_k)_{k=0}^{m-1} = y(x, \tilde{x})$$

of a differentiable program $F(x, \tilde{x})$ decomposes into a sequence of $q = p + m$ differentiable elemental functions $\varphi_j$ evaluated as a single assignment code

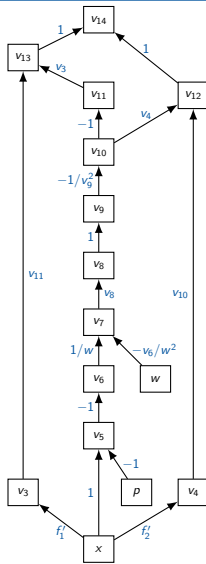$$v_j = \varphi_j(v_k)_{k \prec j} \quad \text{for } j = n, \ldots, n + q - 1$$

over active variables $v_j,\ j = 0, \ldots, n + q - 1$ and where $v_i = x_i$ for $i = 0, \ldots, n - 1$, $y_k = v_{n+p+k}$ for $k = 0, \ldots, m - 1$.

The notation $k \prec j$ $(j \succ k)$ marks $v_k$ as an argument of $\varphi_j$.

```
1   template<typename T>
2   void f(T &x, const T &p, const T &w) {
3       std::vector<T> v(15);
4       v[0]=x;
5       v[1]=p,
6       v[2]=w;
7       v[3]=sin(v[0]);
8       v[4]=cos(v[0]);
9       v[5]=v[0]−v[1];
10      v[6]=−v[5];
11      v[7]=v[6]/v[2];
12      v[8]=exp(v[7]);
13      v[9]=1.0+v[8];
14      v[10]=1.0/v[9];
15      v[11]=1.0−v[10];
16      v[12]=v[3]*v[11];
17      v[13]=v[4]*v[10];
18      v[14]=v[12]+v[13];
19      x=v[14];
20  }
```

Assuming differentiability of all elemental functions the differentiation of

$$v_k = \varphi_k \left( \varphi_j(v_i, (v_\mu)_{i \neq \mu \prec j}), v_i, (v_\nu)_{\{i,j\} \neq \nu \prec k} \right)$$

with respect to $v_i$ yields

$$\frac{dv_k}{dv_i} = \frac{dv_k}{dv_j} \cdot \frac{dv_j}{dv_i} + \frac{\partial v_k}{\partial v_i}$$

where $v_j = \varphi_j(v_i, (v_\mu)_{i \neq \mu \prec j})$.

The corresponding contribution to the directional derivative (tangent) of $v_k$ becomes equal to

$$\frac{dv_k}{dv_i} \left( v_i, (v_\mu)_{i \neq \mu \prec j}, (v_\nu)_{\{i,j\} \neq \nu \prec k} \right) \cdot v_i^{(1)} \ .$$

Example: v_t[7]+=v_t[6]/v[2]; v_t[7]+=−v[6]*v_t[2]/pow(v[2],2);

As an immediate consequence of the chain rule the directional derivative (tangent) of

$$y = (y_k)_{k=0}^{m-1} = y(x, \tilde{x})$$

with respect to $x = (x_i)_{i=0}^{n-1}$ in direction $x^{(1)} = (x_i^{(1)})_{i=0}^{n-1}$ is computed for given $v_i^{(1)} = x_i^{(1)}$ as

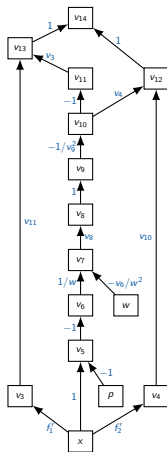$$v_j^{(1)} = v_j^{(1)} + \frac{d\varphi_j}{dv_i}(v_k)_{k \prec j} \cdot v_i^{(1)} \;\; \forall \; i \prec j$$

for $j = n, \dots, n+q-1$ and all $v_j^{(1)}$ equal to zero initially.

The directional derivative is returned through $y^{(1)} = \left(y_k^{(1)}\right)_{k=0}^{m-1}$ where $y_k^{(1)} = v_{n+p+k}^{(1)}$. Tangent arithmetic lends itself to implementation by overloading as a simple augmentation of the primal arithmetic with the computation of directional derivatives.

# Tangent Single Assignment Code
## Example: Sigmoidal Smoothing

```
1   template<typename T>
2   void f_t(T &x, T &x_t,
3            const T &p, const T &p_t, const T &w, const T &w_t) {
4     std::vector<T> v(15), v_t(15,0);
5     v[0]=x; v_t[0]=x_t;
6     v[1]=p; v_t[1]=p_t;
7     v[2]=w; v_t[2]=w_t;
8     v[3]=sin(v[0]); v_t[3]+=cos(v[0])*v_t[0];
9     v[4]=cos(v[0]); v_t[4]+=−sin(v[0])*v_t[0];
10    v[5]=v[0]−v[1]; v_t[5]+=v_t[0]; v_t[5]+=−v_t[1];
11    v[6]=−v[5]; v_t[6]+=−v_t[5];
12    v[7]=v[6]/v[2]; v_t[7]+=v_t[6]/v[2]; v_t[7]+=−v[6]*v_t[2]/pow(v[2],2);
13    v[8]=exp(v[7]); v_t[8]+=v[8]*v_t[7];
14    v[9]=1.0+v[8]; v_t[9]+=v_t[8];
15    v[10]=1.0/v[9]; v_t[10]+=−v_t[9]/pow(v[9],2);
16    v[11]=1.0−v[10]; v_t[11]+=−v_t[10];
17    v[12]=v[3]*v[11]; v_t[12]+=v_t[3]*v[11]; v_t[12]+=v[3]*v_t[11];
18    v[13]=v[4]*v[10]; v_t[13]+=v_t[4]*v[10]; v_t[13]+=v[4]*v_t[10];
19    v[14]=v[12]+v[13]; v_t[14]+=v_t[12]; v_t[14]+=v_t[13];
20    x=v[14]; x_t=v_t[14];
21  }
```

# Tangent Code
## Rules

**TR1** The active data segment must be duplicated. Each active primal variable is matched by its tangent [of same type and shape].

**TR2** Assignment-level tangent code must precede the respective primal assignments. Local tangent single assignment code simplifies differentiation of complex expressions.

**TR3** The tangent flow of control is equal to the primal flow of control. Mind potential non-differentiability due to branching.

**TR4** Calls to primal subprograms must be replaced with calls to the corresponding tangent subprograms. This rule generalizes to polymorphism (overloading, class hierarchies).

**TR5** Drivers are required, e.g, for the accumulation of the gradient.

```
1   template<typename T>
2   void gradient(const std::vector<T> &x, T& y,
3                     std::vector<T> &grad) {
4     size_t n=x.size();
5     std::vector<T> x_t(n,0);
6     for (size_t i=0;i<n;i++) {
7       x_t[i]=1;
8       f_t(x,x_t,y,grad[i]);
9       x_t[i]=0;
10    }
11  }
```

▶ A driver is required to extract the appropriate derivatives from the tangent code in the given context; e.g, the gradient element-wise as directional derivatives in the Cartesian basis directions.

▶ Directional derivatives in other directions may be required.

```
1   void f(float x, float &y) {
2     float z=x*x; y=sin(z);
3   }
4
5   void f_t(float x, float x_t,
6            float &y, float &y_t) {
7     float z_t=2*x*x_t;
8     float z=x*x;
9     y_t=cos(z)*z_t;
10    y=sin(z);
11  }
```

- ▶ The signature is augmented with tangents x_t and y_t for the active arguments x and y.
- ▶ x is passed by value; so is x_t yielding two local variables of type **float**.
- ▶ y is passed by reference; so is y_t which needs to be declared outside of f_t.
- ▶ The local variable z is augmented with its tangent z_t.
- ▶ Both primal assignments are preceded by their (trivial) tangent assignments.

```
1   void f(float x, float &y) {
2       y=sin(x*x);
3   }
4
5   void f_t(float x, float x_t,
6               float &y, float &y_t) {
7       float v_t=2*x*x_t;
8       float v=x*x;
9       y_t=cos(v)*v_t;
10      y=sin(v);
11  }
```

▶ Each (one in this case) primal assignment is decomposed into a single assignment code augmented with its corresponding tangents.

▶ Optimization by *copy propagation* eliminates v_t yielding

```
T v=x*x;
y_t=cos(v)*2*x*x_t;
y=sin(v);
```

► For a given primal implementation of $y = \sin(x^T \cdot x)$ as

```
template<typename T>
void f(const std::vector<T> &x, T &y)
  {
  T xTx=0;
  for (size_t i=0;i<x.size();i++)
    if (i==0)
      xTx=pow(x[i],2);
    else
      xTx+=pow(x[i],2);
  y=sin(xTx);
}
```

we obtain the tangent code on the left.

```
1  template<typename T>
2  void f_t(const std::vector<T> &x, const
      std::vector<T> &x_t, T &y, T &y_t) {
3    T xTx_t=0, xTx=0;
4    for (size_t i=0;i<x.size();i++)
5      if (i==0) {
6        xTx_t=2*x[i]*x_t[i];
7        xTx=pow(x[i],2);
8      } else {
9        xTx_t+=2*x[i]*x_t[i];
10       xTx+=pow(x[i],2);
11     }
12   y_t=cos(xTx)*xTx_t; y=sin(xTx);
13 }
```

► Special care must be taken if the flow of control yields nondifferentiability; e.g, (sigmoidal) smoothing. We focus on differentiable programs.

```
1   void g(float x, float &y) {
2       y=x*x;
3   }
4
5   void g_t(float x, float x_t,
6               float &y, float &y_t) {
7       y_t=2*x*x_t; y=x*x;
8   }
9
10  void f(float x, float &y) {
11      float z; g(x,z); y=sin(z);
12  }
13
14  void f_t(float x, float x_t,
15              float &y, float &y_t) {
16      float z,z_t;
17      g_t(x,x_t,z,z_t);
18      y_t=cos(z)*z_t; y=sin(z);
19  }
```

▶ Calls to primal subprograms need to be replaced by calls to their tangents.

▶ Potentially induced nondifferentiability (e.g, due to recursion) needs to be dealt with. We focus on differentiable programs.

# Outline

In the following we apply the tangent code generation rules to "slightly more real-world" examples, namely the previously introduced implementations of

▶ sigmoidal smoothing to illustrate tangent straight-line code;

▶ Newton's method to illustrate tangent
  ▶ intraprocedural
  ▶ interprocedural
  code.

```
1   template<typename T>
2   void f_t(T &x, T &x_t, const T &p, const T &p_t, const T &w, const T &w_t) {
3       T f1_t=cos(x)*x_t; T f1=sin(x);
4       T f2_t=-sin(x)*x_t; T f2=cos(x);
5       T aux=exp(-(x-p)/w);
6       x_t=-pow(1./(1.+aux),2)*aux*(-x_t/w+p_t/w+w_t*(x-p)/pow(w,2));
7       x=1./(1.+aux);
8       x_t=f1_t*(1-x)-f1*x_t+f2_t*x+f2*x_t;
9       x=f1*(1-x)+f2*x;
10  }
```

# Tangent Intraprocedural Code

Example: Newton's Method

```
1   template<typename T, typename PT>
2   void f_t(T &x, T &x_t, const T &p, const T &p_t, const PT &eps, const unsigned int
        maxit) {
3       unsigned int it=0;
4       T f_t=2*x*x_t−p_t;
5       T f=pow(x,2)−p;
6       do {
7           T dfdx_t=2*x_t;
8           T dfdx=2*x;
9           x_t−=f_t/dfdx−f*dfdx_t/pow(dfdx,2);
10          x−=f/dfdx;
11          f_t=2*x*x_t−p_t;
12          f=pow(x,2)−p;
13          if (++it==maxit) break;
14      } while(fabs(f)>eps);
15  }
```

```
1  template<typename T, typename PT>
2  void f_t(T &x, T& x_t, const PT &p, const PT &p_t, T &r, T &r_t);
3
4  template<typename T, typename PT>
5  void dfdx_t(T &x, T &x_t, const PT &, const PT &, T &drdx, T& drdx_t);
6
7  template<typename T, typename PT>
8  void newton_t(T &x, T &x_t, const T &p, const T &p_t, const PT &eps, const unsigned
        int maxit) {
9      unsigned int it=0;
10     T y,y_t;
11     f_t(x,x_t,p,p_t,y,y_t);
12     do {
13         T dydx,dydx_t;
14         dfdx_t(x,x_t,p,p_t,dydx,dydx_t);
15         x_t-=y_t/dydx-y*dydx_t/pow(dydx,2);
16         x-=y/dydx;
17         f_t(x,x_t,p,p_t,y,y_t);
18         if (++it==maxit) break;
19     } while(fabs(y)>eps);
20 }
```

- association of tangents by address
- vector tangent mode
- assignment-level adjoint code
- lower-precision tangents

# Summary