

Finding gradients by hand

Friday, April 9, 2021

6:37 AM

Overview: how to find gradients?

① Don't find them... use derivative-free methods "DFO"
= Deriv. Free Optim.

These are usually better than finite-diff. methods,
at least if function evaluations are expensive.

Work fine in 1D, OK in small dim., horrible in high dimensions

Usually reserved for last-resort, gradient doesn't exist or hard to obtain

Algos

- old-school stuff like **Nelder-Mead** ≠ simplex (LP) method simplex-reflection
slow, no guarantees
- coordinate descent, pattern search
also slow, not state-of-the-art
- implicit filtering
like finite-diff. but w/o a tiny stepsize h for gradient
- **model-based**
 - **polynomial models** (linear or quadratic)
w/ trust-region
State-of-the-art. Scales to ≈ 1000 dimensions
See Nocedal + Wright or
"Intro to DFO" (SIAM '09) Conn, Scheinberg, Vicente
 - **Gaussian Process model**
w/ **Acquisition Function** "Bayesian Optimization"
trades off exploration / exploitation
Handles discrete variables and noisy functions
(so used often for hyperparameter tuning)
BoTorch software
poor above ≈ 20 dim
- others

- hybrid model/heuristic: **CMA-ES**
- evolutionary search, "zeroth order" methods,
Spall's **SPSA**, Becker's **SSD**
- heuristics: particle-swarm, etc

Often applied to nonconvex problems so bad stationary pts an issue

② Finite difference computation to approximate gradient

Usually not the best solution (ie. poor "long-term" solution) but not to be overlooked, especially for prototyping & back-of-envelope simulations.

Better (ie. cheaper) in small dimensions

Fundamentally low accuracy

Forward diff $f: \mathbb{R}^n \rightarrow \mathbb{R}$ so $\nabla f(x) \in \mathbb{R}^n$ canonical unit basis vector
 $\forall i \in [n], (\nabla f(x))_i \approx \frac{f(x + h e_i) - f(x)}{h}$ so $n+1$ function evaluations

$h \approx 10^{-8}$ rule-of-thumb: too large and Taylor Series is inaccurate
too small and roundoff makes it less accurate

Centered diff is better (keeps about 2/3 of floating pt. digits)
but needs $2n$ function evals

$$\approx \frac{f(x + h e_i) - f(x - h e_i)}{2h}$$

Neither works if f is noisy

③ Compute gradient by hand (or via Mathematica, etc.)

- **Product Rule** $(fg)' = f'g + f \cdot g'$
 - **Leibniz Rule**: functions defined by integrals
 - Functions defined as min/max
- } see later set of notes

Chain Rules

simple 1D $\frac{d}{dx} g \circ f(x) = g'(f(x)) \cdot f'(x)$

multi-variate $\mathbb{R}^n \xrightarrow{f} \mathbb{R}^m \xrightarrow{g} \mathbb{R}^p$
 $\underbrace{\hspace{10em}}_{h = g \circ f}$

Def Jacobian $(Jf)(x)_{ij} = \frac{\partial f_i}{\partial x_j}(x)$ i.e. if $f: \mathbb{R}^n \rightarrow \mathbb{R}^1$
then $(Jf)(x) = \nabla f(x)^T$

then

$(Jh)(x) = (Jg)(f(x)) \cdot (Jf)(x)$ matrix multiplication

special cases: $p=1$ so $Jh = \nabla h^T$

then $\nabla h(x) = (Jf)(x)^T \cdot \nabla g(f(x))$

and if $f(x) := Ax - b$ this simplifies to
 $\nabla h(x) = A^T \cdot \nabla g(Ax - b)$

total derivative version

let $z = f(x, y)$ and $x = g(t)$, $y = h(t)$

then $\frac{dz}{dt} = \frac{\partial f}{\partial x} \cdot \frac{dx}{dt} + \frac{\partial f}{\partial y} \cdot \frac{dy}{dt}$

i.e., implicit differentiation

$F(x, y) = 0$ implicitly defines a function
 $y = f(x)$ (so $F(x, f(x)) = 0$)

under the condition
of implicit fcn then

since $F(x, y) = 0$

$\Rightarrow \frac{dF}{dx}(x, y) = 0$, so chain rule:

$0 = \frac{dF}{dx} = \underbrace{\frac{\partial F}{\partial x}}_{F_x} \cdot \underbrace{\frac{dx}{dx}}_{=1} + \underbrace{\frac{\partial F}{\partial y}}_{F_y} \cdot \frac{dy}{dx}$

i.e. $\frac{dy}{dx} = -\frac{F_x}{F_y}$

Examples

Scalar cases usually straight forward
vector/matrices more work.

See Boyd + Vandenberghe examples such as $\nabla \log \det(X)$

or Peterson + Pederson's "Matrix Cookbook"

or use
Kronecker
product tricks

$$(B^T \otimes A) \text{vec}(X) = \text{vec}(A \cdot X \cdot B)$$

ex: $f(U) = \frac{1}{2} \|U \cdot V^T - B\|_F^2$, $U \in \mathbb{R}^{m \times r}$ is a matrix

idea:

expand

$$f(U + \Delta) = f(U) + \langle \nabla f(U), \Delta \rangle + O(\|\Delta\|^2)$$

$$\langle A, B \rangle = \text{tr}(A^T B)$$

find this term

$$f(U + \Delta) = \frac{1}{2} \|(U + \Delta)V^T - B\|_F^2$$

$$= \frac{1}{2} \|(UV^T - B) + \Delta \cdot V^T\|_F^2$$

$$= \underbrace{\frac{1}{2} \|UV^T - B\|_F^2}_{f(U) \checkmark} + \underbrace{\langle UV^T - B, \Delta \cdot V^T \rangle}_{\text{gradient}} + \underbrace{\frac{1}{2} \|\Delta V^T\|_F^2}_{O(\|\Delta\|^2) \checkmark}$$

$$= \text{tr}((UV^T - B)^T \cdot \Delta \cdot V^T)$$

$$= \text{tr}(V^T (UV^T - B)^T \cdot \Delta) \quad \text{via cyclic property}$$

$$= \text{tr}((UV^T - B) \cdot V)^T \cdot \Delta$$

$$= \langle (UV^T - B) \cdot V, \Delta \rangle$$

so this is the gradient

ie. $\nabla_U \frac{1}{2} \|UV^T - B\|_F^2 = (UV^T - B) \cdot V$

as you'd expect via chain rule, only being careful since matrix multiplication isn't commutative.