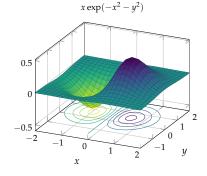
Review: Unconstrained Optimization

Robert Moss Stanford University, Stanford, CA 94305 MOSSR@CS.STANFORD.EDU AA222/CS361

 $\frac{d}{dx}f(g(x)) = \frac{d}{dx}(f \circ g)(x) = \frac{df}{dg}\frac{dg}{dx}$

$$\nabla^2 f(\mathbf{x}) = \mathbf{H}(\mathbf{x}) = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1 \partial x_1} & \frac{\partial^2 f}{\partial x_1 \partial x_2} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2 \partial x_2} \end{bmatrix}$$



¹ A *global optimization method* where *f* is *Lipschitz continuous*.

(2) Derivatives and Gradients.

$$f'(x) \approx \underbrace{\frac{f(x+h) - f(x)}{h}}_{\text{forward difference}} \approx \underbrace{\frac{\text{central difference}}{h}}_{\text{central difference}} \approx \underbrace{\frac{\text{backward difference}}{h}}_{\text{backward difference}}$$

$$f(x) = \text{Re}(f(x+ih)) + h^2 \frac{f''(x)}{2!} - \cdots, \qquad f'(x) = \frac{\text{Im}(f(x+ih))}{h} + O(h^2) \text{ as } h \to 0$$

• *Forward accumulation* will auto-diff f with a single forward pass through computational graph (see *dual numbers:* $a + b\epsilon$ where $\epsilon^2 = 0$). *Reverse accumulation* requires n passes.

(3) Bracketing.

• Fibonacci search, golden section search (Fibonacci approximation), quadratic fit search, Shubert-Piyavskii method,¹ and the bisection method (a root-finding method expanded by Brent-Dekker).

(4) Local Descent.

• Descent direction methods, (approximate) line search (learn step size: $\min_{\alpha} f(\mathbf{x} + \alpha \mathbf{d})$), backtracking line search, Wolfe conditions, and trust regions (or restricted step method).

(5) First-Order Methods.

• Gradient descent, conjugate gradient,³ momentum, Nesterov momentum,⁴ adaptive subgradient method (or Adagrad),⁵ RMSProp,⁶ Adadelta,⁷ adaptive moment estimation (or Adam),⁸ and hypergradient descent (augments a DescentMethod).

(6) Second-Order Methods.

• Newton's method⁹ (quadratic approximation), secant method, ¹⁰ and quasi-Newton methods. ¹¹

⁽⁷⁾Direct Methods.

• Cyclic coordinate search, Powell's method, ¹² Hooke-Jeeves, generalized pattern search, Nelder-Mead simplex method, and divided rectangles (univariate and multivariate DIRECT).

(8) Stochastic Methods.

• Stochastic gradient descent, mesh adaptive direct search, simulated annealing (Metropolis criterion), adaptive simulated annealing, cross-entropy method, ¹³ natural evolutionary strategies, and covariance matrix adaptation evolutionary strategy (CMA-ES: uses a multivariate Gaussian distribution).

(9) Population Methods.

• *Initial population* (uniform, normal, Cauchy), *genetic algorithms* (undergo *crossover* and *mutation*), *differential evolution*, *particle swarm optimization*, *firefly algorithm*, *cuckoo search*, *hybrid methods* (local search using descent methods: *Lamarchian learning* and *Baldwinian learning*).

² Start big then backtrack.

³ Overcomes narrow valley issues of *gradient descent*. Its directions are *mutually conjugate* with respect to **A**. Approximations of β are *Fletcher-Reeves* and *Polak-Ribière*.

⁴ Reduce overshooting at bottom.

⁵ Dulls high gradients, increases influence of infreq. updated params.

⁶ Extends *Adagrad* to avoid effects of a monotonically decreasing learning rate.

 7 Overcomes *Adagrad's* monotonically decreasing learning rate by eliminating α in favor of:

$$x' = x - \frac{\text{RMS}(\Delta x)}{\epsilon + \text{RMS}(g)}g$$

 8 Adapts learning rate to each parameter, biases then corrects decay in momentum ${\bf v}$ and sq. gradient ${\bf s}.$

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

¹⁰ Unlike *Newton's method*, estimates f'', only requires f':

$$f''(x) \approx \frac{f'(x) - f'(x^{(k-1)})}{x - x^{(k-1)}}$$

 11 Inverse Hessian \mathbf{H}^{-1} approximation: *DFP, BFGS, L-BFGS* approximate line search.

12 Non-orthogonal directions.

¹³ Sample, select, fit, repeat.