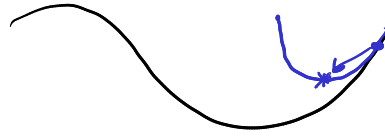


Last time

Newton's Method

Quasi-Newton

Approximate H^{-1}



Conjugate Gradient and Momentum

Line Search

Sat. Wolfe \Rightarrow Global Convergence
(Thm 3.2)

Direction Affects Convergence
Rate

Trust Region

Momentum

- Not in N+W (2006)
- Training of Neural Networks
- Availability of Auto-diff

Corolla

Computational
Cost
sensitive to n

Cadillac (Luxury)

Directions for Line Search

Gradient (Steepest)

Conjugate Gradient

L-BFGS

Quasi Newton

BFGS

Newton

Convergence Rate

Linear

No theorem slower
in N.W.
(for quadratic, LCG
converges in n steps)

Super linear

Quadratic

Conjugate Gradient

$$f(\vec{x}) = \frac{1}{2} \vec{x}^T A \vec{x} + \vec{b}^T \vec{x} + c$$

$$\nabla f(\vec{x}) = A \vec{x} + \vec{b} = \vec{0}$$

$$\vec{x}^* = -A^{-1} \vec{b}$$

Newton: Invert

Conjugate Direction

Def: \vec{d}^i and \vec{d}^j are mutually conjugate w.r.t. A if

$$\vec{d}^i^T A \vec{d}^j = 0$$

(not generally orthogonal)

Linear Conjugate Gradient (LCG)

start with $\vec{d}^1 = -\vec{g}^1$
 loop
 $\vec{x}^{k+1} \leftarrow \vec{x}^k + \alpha^k \vec{d}^k$
 $\vec{d}^{k+1} \leftarrow -\vec{g}^{k+1} + \beta^{k+1} \vec{d}^k$

α^k is minimizer along \vec{d}^k

$$\alpha^k = - \frac{\vec{d}^k T (A \vec{x}^k + \vec{b})}{\vec{d}^k T A \vec{d}^k}$$

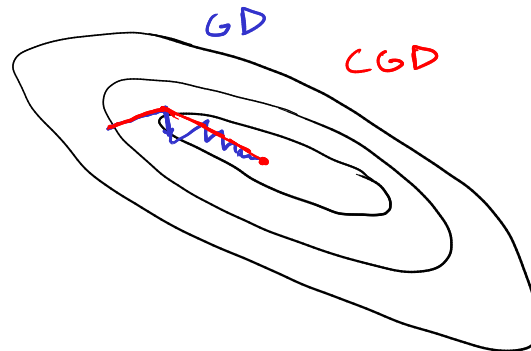
$$\vec{d}^k T A \vec{d}^{k-1} = 0$$

$$(-\vec{g}^k + \beta^k \vec{d}^{k-1}) T A \vec{d}^{k-1} = 0$$

$$-\vec{g}^k T A \vec{d}^{k-1} + \beta^k \vec{d}^{k-1 T} A \vec{d}^{k-1} = 0$$

$$\beta^k = \frac{\vec{g}^k T A \vec{d}^{k-1}}{\vec{d}^{k-1 T} A \vec{d}^{k-1}}$$

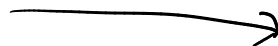
LCG finds minimum of f in n steps
 (N+W Thm 5.1)



Large Nonlinear Problems: Unknown A (or expensive)
 $\hookrightarrow \alpha? \beta?$
 - Approximate α with line search
 - Approximate β (below)

Fletcher Reeves

$$\beta^k = \frac{\vec{g}^k T \vec{g}^k}{\vec{g}^{k-1 T} \vec{g}^{k-1}}$$



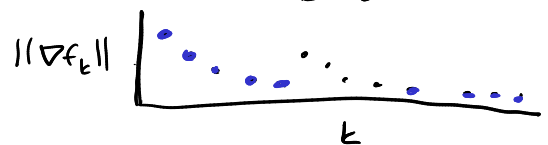
N.W. Theorem 5.7

Strong Wolfe line search

$$\Rightarrow \liminf_{k \rightarrow \infty} \|\nabla f_k\| = 0$$

Polak-Ribière

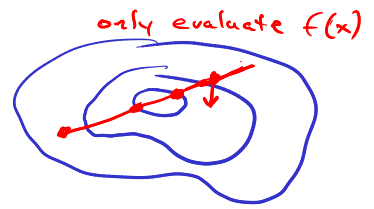
$$\beta^k = \max\left(\frac{\vec{g}^k T (\vec{g}^k - \vec{g}^{k-1})}{\vec{g}^{k-1 T} \vec{g}^{k-1}}, 0\right)$$



Relationship between CG and Quasi-Newton
 p. 180 of N+W

"Memoriless" BFGS (at each step, set $\vec{Q}^k = I$)
 then equivalent to Hestenes-Stiefel CG

Break! If you have R.M. Autodiff why would you prefer fixed step size gradient descent over ~~line search~~ backtracking line search



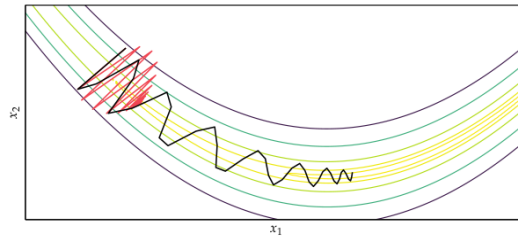
Momentum / Adaptive Scaling

"Vanilla"

Momentum

$$\begin{aligned}\vec{v}^{k+1} &\leftarrow \beta \vec{v} - \alpha \vec{g}^k \\ \vec{x}^{k+1} &\leftarrow \vec{x}^k + \vec{v}^{k+1}\end{aligned}$$

Tendency to overshoot



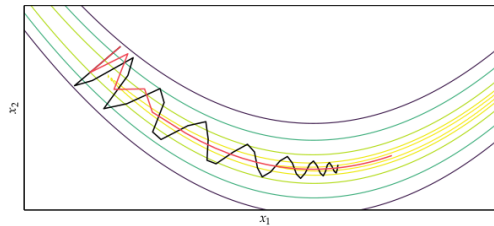
— gradient descent
— momentum

Figure 5.5. Gradient descent and the momentum method compared on the Rosenbrock function with $b = 100$; see appendix B.6.

Momentum

Nesterov Momentum

$$\begin{aligned}\vec{v}^{k+1} &\leftarrow \beta \vec{v}^k - \alpha \nabla f(\vec{x}^k + \beta \vec{v}^k) \\ \vec{x}^{k+1} &\leftarrow \vec{x}^k + \vec{v}^{k+1}\end{aligned}$$



— momentum
— Nesterov momentum

Figure 5.6. The momentum and Nesterov momentum methods compared on the Rosenbrock function with $b = 100$; see appendix B.6.

Adagrad

$$x_i^{k+1} \leftarrow x_i^k - \frac{\alpha}{\epsilon + \sqrt{s_i^k}} g_i^k$$

$$s_i^k = \sum_{j=1}^k (g_i^j)^2$$

monotonic decrease of learning rate

monotonically increases

element-wise

RMS Prop

$$\hat{s}_i^k = \gamma \hat{s}_i^{k-1} + (1-\gamma)(g_i^k)^2$$

$$\gamma = 0.9$$

$$\text{RMS}(g_i) \equiv \sqrt{\hat{s}_i}$$

Step size adaptation

Ada delta

$$\vec{x}_i^{k+1} \leftarrow \vec{x}_i^k - \frac{\text{RMS}(\Delta x_i)}{\epsilon + \text{RMS}(g_i)} g_i^k$$

No learning rate

Ada m

$$\vec{v}^{k+1} \leftarrow \gamma_v \vec{v}^k + (1-\gamma_v) \vec{g}^k$$

$$\vec{s}^{k+1} \leftarrow \gamma_s \vec{s}^k + (1-\gamma_s) (\vec{g}^k \odot \vec{g}^k)$$

$$\hat{\vec{v}}^{k+1} \leftarrow \frac{\vec{v}^{k+1}}{1-\gamma_v}$$

$$\hat{\vec{s}}^{k+1} \leftarrow \frac{\vec{s}^{k+1}}{1-\gamma_s}$$

$$\vec{x}^{k+1} = \vec{x}^k - \frac{\alpha \hat{\vec{v}}^{k+1}}{\epsilon + \sqrt{\hat{\vec{s}}^{k+1}}}$$