Classical Gradient Methods

- Note simultaneous course at AMSI (math) summer school: Nonlin. Optimization Methods (see <u>http://wwwmaths.anu.edu.au/events/amsiss05/</u>)
- Recommended textbook (Springer Verlag, 1999): Nocedal & Wright, Numerical Optimization
- Here: just quick overview, unconstrained only
- But will consider large, nonlinear problems

Function Optimization

- Goal: given (diff able) function $f: \mathbb{R}^{n} \to \mathbb{R}$ find minimum $\overline{\mathscr{O}}^{*} = \arg\min f(\overline{\mathscr{O}})$
- Gradient methods find only local minimum
- For convex functions, local min. = global min.
 In machine learning,
- Fn. is defined over data: $f(\vec{w}) = B_{\vec{w}}[f(\vec{w}; \vec{x})]$
- May comprise loss and regularization terms

Methods by Gradient Order

- Oth order (direct, gradient-free) methods use only the function values themselves
- 1st order gradient methods additionally use function's gradient $\vec{g} = \vec{g}(\vec{w}) =$
- 2nd order gradient methods also use the function's Hessian

$$H=H(ec{w})=rac{\partial^2 f(ec{w})}{\partial ec{w}\,\partial ec{w}^1}$$

Direct (Gradient-Free) Methods

Many distinct algorithms:

- Simulated annealing, Monte Carlo optim.
- Perturbation methods, SPSA, Tabu search
- Genetic algorithms, evolutionary strategies, ant colony optimization, ...

Differ in many implementation details but all share a common approach.

Prototypical Direct Method

Randomly initialize pool W of candidates Repeat until converged:

pick parent(s) w_i from generate child(ren) $w_i' = perturb(w_i)$ compare child to parent (or entire pool): $\Delta = f(w_i') - f(w_i)$ if $\Delta < 0$ accept w_i' into W (may replace w_i) else if global optimization:

accept w_i' into W with probability P($e^{-\Delta}$)

Direct Methods: Advantages

- No need to derive or compute gradients
 - Can solve discrete/combinatorial problems
 - Can address even non-formalized problems
- Can find (non-convex fn.'s) global optimum
- Highly and easily parallelizable
- Very fast iteration when perturbation and evaluation are both incremental, *i.e.* O(1)

Direct Methods: Disadvantages

- No sense of appropriate direction or size of step to take (perturbation is random)
 - Some algorithms try to fix this heuristically
- Takes too many iterations to converge
- Global optim. requires knowing acceptance of inferior candidates ⇒ slower still
- No strong mathematical underpinnings
 ⇒ jungle of ad-hoc heuristics

Gradient Descent

- Perturbs parameter vector in steepest downhill direction (= neg. gradient): \$\vec{w}_{k+1} = \vec{w}_k \eta \vec{g}_k\$
- Step size η can be set
 - to small positive constant: simple gradient descent
 - by line minimization: steepest descent
 - adaptively (more on this later)

Advantage:

• Cheap to compute: iteration typically just O(n)

Gradient Descent: Disadvantages

- Line minimization may be expensive
- Convergence slow for ill-conditioned problems: #iterations \geq condition# $\kappa = \frac{\lambda_{max}}{\lambda_{max}}$ of Hessian



Newton's Method

• Local quadratic model $f(\vec{w}) = \frac{1}{2}(\vec{w} - \vec{w}^*)^T H(\vec{w} - \vec{w}^*)$ has gradient $\vec{g}(\vec{w}) = H(\vec{w} - \vec{w}^*)$ therefore let $\vec{w}_{i+1} = \vec{w}_i - H_i^{-1} \vec{g}_i$

Newton's Method

Big advantage:

• Jumps directly to minimum of quadratic bowl (regardless of ill-conditioning)

Disadvantages:

- Hessian expensive to invert: nearly $O(n^3)$
- Hessian must be positive definite: $\lambda_{min} > 0$
- May make huge, uncontrolled steps



Gauss-Newton Approximation

- Let $f = l \circ \vec{m}$, $\vec{m} : \mathbb{R}^{n} \to \mathbb{R}^{p}$ = model, l = lossThen $H_{f} = \underbrace{J_{\vec{m}}^{T} H_{l} J_{\vec{m}}}_{\mathcal{H}_{l} = 1} + \sum_{i=1}^{p} H_{m_{i}}(J_{l})_{i}$ Gauss-Newton: G_{f} Jacobian: • $H_{l} \ge 0 \Rightarrow G_{f} \ge 0$ $(J_{\vec{m}})_{ij} = \frac{\partial m_{i}(\vec{w})}{\partial m_{i}}$
- At minimum, $G_f = H_f$
- For sum-squared loss: $H_l = I$ pseudo-inverse and $G_{\vec{j}}^{-1}\vec{j} = (J_{\vec{m}}^T J_{\vec{m}})^{-1} J_{\vec{m}}^T J_{\vec{l}}^T = J_{\vec{m}}^+ J_{\vec{l}}^T$

Levenberg-Marquardt

$$\vec{w}_{t+1} = \vec{w}_t - (G_t + \lambda \operatorname{diag}(G_t))^{-1} \vec{g}_t$$

- G_t is Gauss-Newton approximation to H_t (guaranteed positive semi-definite)
- $\lambda \ge 0$ adaptively controlled, limits step to an elliptical model-trust region
- Fixes Newton's stability issues, but still $O(n^3)$

Quasi-Newton: BFGS

- Iteratively updates estimate B of H^{-1}
- Guarantees $B^T = B$ and B > 0
- Reduces complexity to $O(n^2)$ per iteration
- Requires line minimization (direction $B_t \vec{g}_t$)
- Update formula:

$$B_{t+1} = B_t + \frac{\bigtriangleup \vec{w} \bigtriangleup \vec{w}^T}{\bigtriangleup \vec{w}^T \bigtriangleup \vec{y}} - \frac{B_t \bigtriangleup \vec{y} \bigtriangleup \vec{j}^T B_t^T}{\bigtriangleup \vec{y}^T B_t \bigtriangleup \vec{y}} + \vec{u} \bigtriangleup \vec{j}^T B_t \bigtriangleup \vec{y} \vec{u}^T$$
where
$$\vec{u} \equiv \frac{\bigtriangleup \vec{w}}{\bigtriangleup \vec{w}^T \bigtriangleup \vec{y}} - \frac{B_t \bigtriangleup \vec{y}}{\bigtriangleup \vec{y}^T B_t \bigtriangleup \vec{y}}$$
16

Conjugate Gradient

 $\vec{w}_{i+1} = \vec{w}_i + \alpha \vec{v}_i$; α set by line minimization

Search directions $\vec{v}_0 = -\vec{j}_0; \vec{v}_{i+1} = \beta \vec{v}_i - \vec{j}_i$ are conjugate:

$$i \neq j \Rightarrow \vec{v}_i^T H \vec{v}_j = 0$$

$$eta = rac{ec{g}_t^T(ec{g}_t - ec{g}_{t-1})}{ec{v}_t^T(ec{g}_t - ec{g}_{t-1})}$$

(= orthogonal in local Mahalonobis metric) (Hestenes-Stiefel, 1952) (a.k.a. Beale-Sørenson)

NB: other formulae for β (Polak-Ribiere, Fletcher-Reeves) equivalent for quadratic but inferior for nonlinear fn.s!

Conjugate Gradient: Properties

- No matrices \Rightarrow each iteration costs only O(n)
- Minimizes quadratic fn. exactly in *n* iterations
- Restart every *n* iterations for nonlinear fn.s
- Optimal progress after k < n iterations

An incremental 2nd-order method! Revolutionary.

• Drives nearly all large-scale optimization today.