

Optimization Algorithms

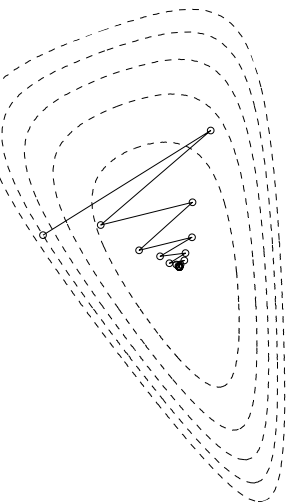
Unconstraint Optimization Basics

Descent direction & stepsize, plain gradient descent, stepsize adaptation & backtracking line search, trust region, steepest descent, Newton, Gauss-Newton, Quasi-Newton, BFGS, conjugate gradient, exotic: Rprop

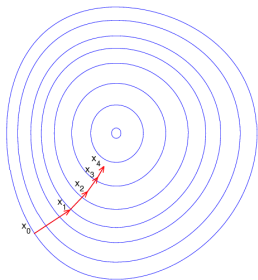
Marc Toussaint

Technical University of Berlin

Winter 2020/21



Gradient descent



- Objective function: $f : \mathbb{R}^n \rightarrow \mathbb{R}$
Gradient vector: $\nabla f(x) = \left[\frac{\partial}{\partial x} f(x) \right]^T \in \mathbb{R}^n$
- Problem:

$$\min_x f(x)$$

where we can evaluate $f(x)$ and $\nabla f(x)$ for any $x \in \mathbb{R}^n$

- Plain gradient descent: iterative steps in the direction $-\nabla f(x)$.

Input: initial $x \in \mathbb{R}^n$, function $\nabla f(x)$, stepsize α , tolerance θ

Output: x

1: **repeat**

2: $x \leftarrow x - \alpha \nabla f(x)$

3: **until** $|\Delta x| < \theta$ [perhaps for 10 iterations in sequence]

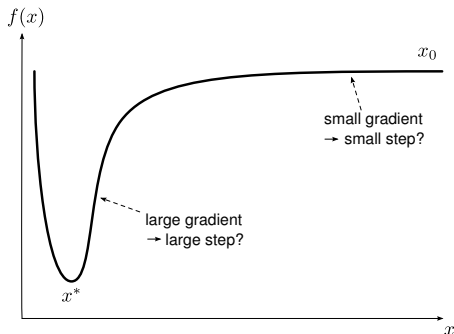
- Plain gradient descent is really not efficient
- Two core issues of unconstrained optimization:

A. Step size

B. Descent direction

Stepsize

- Making steps proportional to $\nabla f(x)$?



- We need methods that
 - robustly adapt stepsize
 - exploit convexity, if known
 - perhaps be independent of $|\nabla f(x)|$ (e.g. if non-convex as above)

Stepsize Adaptation: Backtracking Line Search

Input: initial $x \in \mathbb{R}^n$, functions $f(x)$ and $\nabla f(x)$, tolerance θ , parameters
(defaults: $\varrho_{\alpha}^+ = 1.2$, $\varrho_{\alpha}^- = 0.5$, $\delta_{\max} = \infty$, $\varrho_{\text{ls}} = 0.01$)

- 1: initialize stepsize $\alpha = 1$
- 2: **repeat**
- 3: $\delta \leftarrow -\frac{\nabla f(x)}{|\nabla f(x)|}$ // (alternative: $\delta = -\nabla f(x)$)
- 4: **while** $f(x + \alpha\delta) > f(x) + \varrho_{\text{ls}} \nabla f(x)^\top (\alpha\delta)$ **do** // **line search**
- 5: $\alpha \leftarrow \varrho_{\alpha}^- \alpha$ // decrease stepsize
- 6: **end while**
- 7: $x \leftarrow x + \alpha\delta$
- 8: $\alpha \leftarrow \min\{\varrho_{\alpha}^+ \alpha, \delta_{\max}\}$ // increase stepsize
- 9: **until** $|\alpha\delta| < \theta$ // perhaps for 10 iterations in sequence

- α determines the absolute stepsize
- Guaranteed monotonicity (by construction)
("Typically" ensures convergence to locally convex minima; see later)

Backtracking line search

- Line search in general denotes the problem

$$\min_{\alpha \geq 0} f(x + \alpha\delta)$$

for some step direction δ .

- The most common line search is **backtracking**, which decreases α as long as

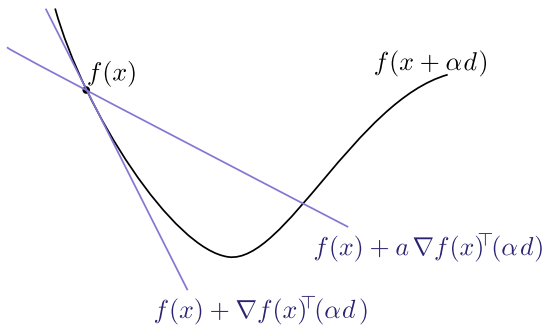
$$f(x + \alpha\delta) > f(x) + \varrho_{\text{ls}} \nabla f(x)^\top (\alpha\delta)$$

ϱ_α^- describes the stepsize decrement in case of a rejected step

ϱ_{ls} describes a minimum desired decrease in $f(x)$

- Boyd et al: typically $\varrho_{\text{ls}} \in [0.01, 0.3]$ and $\varrho_\alpha^- \in [0.1, 0.8]$

Backtracking line search



Wolfe Conditions

- The 1st Wolfe condition (“sufficient decrease condition”)

$$f(x + \alpha\delta) \leq f(x) + \rho_{1s} \nabla f(x)^\top (\alpha\delta)$$

requires a decrease of f at least ρ_{1s} -times “as expected”

- The 2nd (stronger) Wolfe condition (“curvature condition”)

$$|\nabla f(x + \alpha\delta)^\top \delta| \leq \rho_{2s} |\nabla f(x)^\top \delta|$$

implies a requires an decrease of the slope by a factor ρ_{2s} .

$\rho_{2s} \in (\rho_{1s}, \frac{1}{2})$ (for conjugate gradient)

- See Nocedal et al., Section 3.1 & 3.2 for more general proofs of convergence of any method that ensures the Wolfe conditions after each line search

Convergence for (locally) convex functions

- **Theorem** (Exponential convergence on convex functions)
 - Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be an objective function
 - with eigenvalues λ of the Hessian $\nabla^2 f(x)$ bounded by $m < \lambda < M$, with $m > 0, \forall x \in \mathbb{R}^n$
 - then gradient descent with backtracking line search converges exponentially with convergence rate $(1 - 2\frac{m}{M} \varrho_{\text{ls}} \varrho_{\alpha}^-)$.
- The exercises guide through the proof

Newton Methods & Descent Direction

Steepest Descent Direction

- The gradient $-\nabla f(x)$ is sometimes called *steepest descent direction*

Is it really?

Steepest Descent Direction

- The gradient $-\nabla f(x)$ is sometimes called *steepest descent direction*

Is it really?

- Here is a possible definition:

*The steepest descent direction is the one where, **when you make a step of length 1**, you get the largest decrease of f in its linear approximation.*

$$\operatorname{argmin}_{\delta} \nabla f(x)^{\top} \delta \quad \text{s.t. } \|\delta\| = 1$$

Steepest Descent Direction

- But the norm $\|\delta\|^2 = \delta^\top A \delta$ depends on the metric A !

Let $A = B^\top B$ (Cholesky decomposition) and $z = B\delta$

$$\begin{aligned}\delta^* &= \underset{\delta}{\operatorname{argmin}} \nabla f^\top \delta && \text{s.t. } \delta^\top A \delta = 1 \\ &= B^{-1} \underset{z}{\operatorname{argmin}} (B^{-1} z)^\top \nabla f && \text{s.t. } z^\top z = 1 \\ &= B^{-1} \underset{z}{\operatorname{argmin}} z^\top B^{-\top} \nabla f && \text{s.t. } z^\top z = 1 \\ &= B^{-1} [-B^{-\top} \nabla f] = -A^{-1} \nabla f\end{aligned}$$

- The steepest descent direction is $\delta = -A^{-1} \nabla f$

Behavior under linear coordinate transformations

- Let B be a matrix that describes a linear transformation in coordinates
- A coordinate vector x transforms as $z = Bx$
- The gradient vector $\nabla_x f(x)$ transforms as $\nabla_z f(z) = B^{-\top} \nabla_x f(x)$
- The metric A transforms as $A_z = B^{-\top} A_x B^{-1}$
- The steepest descent transforms as $A_z^{-1} \nabla_z f(z) = B A_x^{-1} \nabla_x f(x)$

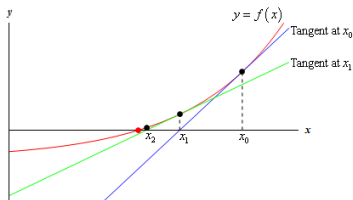
The steepest descent transforms like a normal coordinate vector (covariant)

(more details in the *Maths script*)

Newton Step

Newton Step

- For finding roots (zero points) of $f(x)$



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

- For finding optima of $f(x)$ in 1D (which are roots of $f'(x)$):

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

- For finding optima in higher dimensions $x \in \mathbb{R}^n$:

$$x \leftarrow x - \nabla^2 f(x)^{-1} \nabla f(x)$$

Newton Step

- Assume we have access to the symmetric **Hessian**

$$\nabla^2 f(x) = \begin{pmatrix} \frac{\partial^2}{\partial x_1 \partial x_1} f(x) & \frac{\partial^2}{\partial x_1 \partial x_2} f(x) & \cdots & \frac{\partial^2}{\partial x_1 \partial x_n} f(x) \\ \frac{\partial^2}{\partial x_1 \partial x_2} f(x) & & & \vdots \\ \vdots & & & \vdots \\ \frac{\partial^2}{\partial x_n \partial x_1} f(x) & \cdots & \cdots & \frac{\partial^2}{\partial x_n \partial x_n} f(x) \end{pmatrix} \in \mathbb{R}^{n \times n}$$

- which defines the Taylor expansion:

$$f(x + \delta) \approx f(x) + \nabla f(x)^\top \delta + \frac{1}{2} \delta^\top \nabla^2 f(x) \delta$$

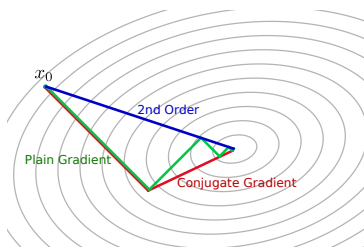
Note: $\nabla^2 f(x)$ acts like a **metric** for δ

Notes on the Newton Step

- If f is a 2nd-order polynomial, the Newton step jumps to the optimum in just one step.
- *Unlike the gradient magnitude* $|\nabla f(x)|$, the magnitude of the Newton step δ is meaningful and scale invariant!
 - If you'd rescale f (trade cents by Euros), δ is unchanged
 - $|\delta|$ is the distance to the optimum of the 2nd-order Taylor.
- *Unlike the gradient* $\nabla f(x)$, the Newton step δ is truly a vector!
 - The Newton step is invariant under coordinate transformations; the coordinates of δ transforms contra-variant, as it is supposed to for vector coordinates
 - The proof is exactly the same as for the steepest descent with a non-Euclidean metric – the Hessian acts as a metric

Why 2nd order information is better

- Better direction:



- Better stepsize:

- A full Newton step jumps directly to the minimum of the local squared approx.
- Robust Newton methods combine this with line search and damping (Levenberg-Marquardt)

Basic Newton method

Input: initial $x \in \mathbb{R}^n$, functions $f(x)$, $\nabla f(x)$, $\nabla^2 f(x)$, tolerance θ , parameters (defaults: $\varrho_\alpha^+ = 1.2$, $\varrho_\alpha^- = 0.5$, $\varrho_{\text{ls}} = 0.01$, λ)

- 1: initialize stepsize $\alpha = 1$, fixed damping λ
- 2: **repeat**
- 3: compute δ to solve $(\nabla^2 f(x) + \lambda \mathbf{I}) \delta = -\nabla f(x)$
- 4: **while** $f(x + \alpha\delta) > f(x) + \varrho_{\text{ls}} \nabla f(x)^\top (\alpha\delta)$ **do** *// line search*
- 5: $\alpha \leftarrow \varrho_\alpha^- \alpha$ *// decrease stepsize*
- 6: **end while**
- 7: $x \leftarrow x + \alpha\delta$ *// step is accepted*
- 8: $\alpha \leftarrow \min\{\varrho_\alpha^+ \alpha, 1\}$ *// increase stepsize*
- 9: **until** $\|\alpha\delta\|_\infty < \theta$

- Notes:

- Line 3 computes the Newton step $\delta = -\nabla^2 f(x)^{-1} \nabla f(x)$,
e.g. using a special Lapack routine `dposv` to solve $Ax = b$ (using Cholesky)

Damping, λ , Trust-Region, Levenberg-Marquardt

- δ solves the problem

$$\min_{\delta} \left[\nabla f(x)^\top \delta + \frac{1}{2} \delta^\top \nabla^2 f(x) \delta + \frac{1}{2} \lambda \delta^2 \right].$$

- λ introduces a squared penalty for large steps

Damping, λ , Trust-Region, Levenberg-Marquardt

- δ solves the problem

$$\min_{\delta} \left[\nabla f(x)^\top \delta + \frac{1}{2} \delta^\top \nabla^2 f(x) \delta + \frac{1}{2} \lambda \delta^2 \right].$$

- λ introduces a squared penalty for large steps

- **Trust region method:**

$$\min_{\delta} \left[\nabla f(x)^\top \delta + \frac{1}{2} \delta^\top \nabla^2 f(x) \delta \right] \quad \text{s.t.} \quad \delta^2 \leq \beta$$

- β defines the *trust region*

Damping, λ , Trust-Region, Levenberg-Marquardt

- δ solves the problem

$$\min_{\delta} \left[\nabla f(x)^\top \delta + \frac{1}{2} \delta^\top \nabla^2 f(x) \delta + \frac{1}{2} \lambda \delta^2 \right].$$

– λ introduces a squared penalty for large steps

- **Trust region method:**

$$\min_{\delta} \left[\nabla f(x)^\top \delta + \frac{1}{2} \delta^\top \nabla^2 f(x) \delta \right] \quad \text{s.t.} \quad \delta^2 \leq \beta$$

– β defines the *trust region*

- Solving this using Lagrange parameters (as we will learn it later):

$$\begin{aligned} L(\delta, \lambda) &= \nabla f(x)^\top \delta + \frac{1}{2} \delta^\top \nabla^2 f(x) \delta + \lambda(\delta^2 - \beta) \\ \nabla_{\delta} L(\delta, \lambda) &= \nabla f(x)^\top + \delta^\top (\nabla^2 f(x) + 2\lambda \mathbf{I}) \end{aligned}$$

gives the step $\delta = -(\nabla^2 f(x) + 2\lambda \mathbf{I})^{-1} \nabla f(x)$, with λ the **dual variable**

- For $\lambda \rightarrow \infty$, δ becomes aligned with $-\nabla f(x)$ (but $|\delta| \rightarrow 0$)

Newton method with non-pos-def fallback

```
1: initialize stepsize  $\alpha = 1$ 
2: repeat
3:   try to compute  $\delta$  to solve  $(\nabla^2 f(x) + \lambda \mathbf{I}) \delta = -\nabla f(x)$ 
4:   if  $\nabla f(x)^\top \delta > 0$  (non-descent) or fails (ill-def. linear system) then
5:      $\delta \leftarrow -\frac{\nabla f(x)}{|\nabla f(x)|}$  // (gradient direction)
6:     (Or: choose  $\lambda > [-\text{minimal eigenvalue of } \nabla^2 f(x)]^+$  and repeat)
7:   end if
8:   while  $f(x + \alpha\delta) > f(x) + \varrho_{\text{ls}} \nabla f(x)^\top (\alpha\delta)$  do // line search
9:      $\alpha \leftarrow \varrho_{\alpha}^- \alpha$  // decrease stepsize
10:    optionally:  $\lambda \leftarrow \varrho_{\lambda}^+ \lambda$  and recompute  $d$  // increase damping
11:  end while
12:   $x \leftarrow x + \alpha\delta$  // step is accepted
13:   $\alpha \leftarrow \min\{\varrho_{\alpha}^+ \alpha, 1\}$  // increase stepsize
14:  optionally:  $\lambda \leftarrow \varrho_{\lambda}^- \lambda$  // decrease damping
15: until  $\|\alpha\delta\|_{\infty} < \theta$  repeatedly
```

- See comments in exercise e02
- Demo ...

- In the remainder: Extensions of the Newton approach:
 - Gauss-Newton
 - Quasi-Newton
 - BFGS, (L)BFGS
 - Conjugate Gradient

- And a crazy method: Rprop

Gauss-Newton method

- Consider a **sum-of-squares** problem:

$$\min_x f(x) \quad \text{where } f(x) = \phi(x)^\top \phi(x) = \sum_{i=1}^d \phi_i(x)^2$$

with **features** $\phi(x) \in \mathbb{R}^d$, and we can evaluate $\phi(x)$, $\frac{\partial}{\partial x} \phi(x)$ for any $x \in \mathbb{R}^n$

- $\phi(x) \in \mathbb{R}^d$ is a vector; each entry contributes a squared cost term to $f(x)$
- $\frac{\partial}{\partial x} \phi(x)$ is the **Jacobian** ($d \times n$ -matrix)

$$\frac{\partial}{\partial x} \phi(x) = \begin{pmatrix} \frac{\partial}{\partial x_1} \phi_1(x) & \frac{\partial}{\partial x_2} \phi_1(x) & \cdots & \frac{\partial}{\partial x_n} \phi_1(x) \\ \frac{\partial}{\partial x_1} \phi_2(x) & & & \vdots \\ \vdots & & & \vdots \\ \frac{\partial}{\partial x_1} \phi_d(x) & \cdots & \cdots & \frac{\partial}{\partial x_n} \phi_d(x) \end{pmatrix} \in \mathbb{R}^{d \times n}$$

with 1st-order Taylor expansion $\phi(x + \delta) = \phi(x) + \frac{\partial}{\partial x} \phi(x) \delta$

Gauss-Newton method

- The gradient and Hessian of $f(x)$ are

$$f(x) = \phi(x)^\top \phi(x)$$

$$\nabla f(x) = 2 \frac{\partial}{\partial x} \phi(x)^\top \phi(x) \quad (\text{recall } \nabla f(x) \equiv \frac{\partial}{\partial x} f(x)^\top)$$

$$\nabla^2 f(x) = 2 \frac{\partial}{\partial x} \phi(x)^\top \frac{\partial}{\partial x} \phi(x) + 2\phi(x)^\top \nabla^2 \phi(x)$$

- *The Gauss-Newton method is the Newton method for $f(x) = \phi(x)^\top \phi(x)$ while approximating $\nabla^2 \phi(x) \approx 0$*

In the Newton algorithm, replace line 3 by solving

$$\left(2 \frac{\partial}{\partial x} \phi(x)^\top \frac{\partial}{\partial x} \phi(x) + \lambda \mathbf{I}\right) \delta = -2 \frac{\partial}{\partial x} \phi(x)^\top \phi(x)$$

Gauss-Newton method

- The approximate Hessian $H = 2 \frac{\partial}{\partial x} \phi(x)^\top \frac{\partial}{\partial x} \phi(x)$ is **always semi-pos-def!**

Gauss-Newton method

- The approximate Hessian $H = 2 \frac{\partial}{\partial x} \phi(x)^\top \frac{\partial}{\partial x} \phi(x)$ is **always semi-pos-def!**
- H is a sum of rank-1 matrices:

$$H = 2 \sum_{i=1}^d \nabla \phi_i(x)^\top \nabla \phi_i(x)$$

(which implies semi-pos-def)

Gauss-Newton method

- The approximate Hessian $H = 2 \frac{\partial}{\partial x} \phi(x)^\top \frac{\partial}{\partial x} \phi(x)$ is **always semi-pos-def!**
- H is a sum of rank-1 matrices:

$$H = 2 \sum_{i=1}^d \nabla \phi_i(x)^\top \nabla \phi_i(x)$$

(which implies semi-pos-def)

- Computing H requires only first-order derivatives of features ϕ , no computationally expensive Hessians

Gauss-Newton method

- The approximate Hessian $H = 2 \frac{\partial}{\partial x} \phi(x)^\top \frac{\partial}{\partial x} \phi(x)$ is **always semi-pos-def!**
- H is a sum of rank-1 matrices:

$$H = 2 \sum_{i=1}^d \nabla \phi_i(x)^\top \nabla \phi_i(x)$$

(which implies semi-pos-def)

- Computing H requires only first-order derivatives of features ϕ , no computationally expensive Hessians
- H can be interpreted as pullback of Euclidean norm $\phi^\top \phi$ in feature space. As it is x -dependent, this is a non-constant metric in x -space – it defines a *Riemannian* metric. (See math notes.)

Robotics example

- Sum-of-squares problems appear at many places in AI: “minimize squared errors”
- Basic robotics example: Trajectory optimization

Let $q : \{1, \dots, T\} \mapsto \mathbb{R}^n$ be a discretized path in \mathbb{R}^n

$$\min_q \sum_{t=1}^T (q_t + q_{t-2} - 2q_{t-1})^2 + \phi(q_T)^2$$

where q_0, q_{-1} are given, and $\phi(q_T)$ are some features of the end configuration q_T

Quasi-Newton methods

Quasi-Newton methods

- Assume we *cannot* evaluate $\nabla^2 f(x)$.

Can we still use 2nd order methods?

- Yes: We can approximate $\nabla^2 f(x)$ from the data $\{(x_i, \nabla f(x_i))\}_{i=1}^k$ of previous iterations

(General view: We can *learn* from the data $\{(x_i, \nabla f(x_i))\}_{i=1}^k \rightsquigarrow$ e.g., Bayesian optimization or model-based for blackbox optimization)

Basic example

- We've seen already two data points $(x_1, \nabla f(x_1))$ and $(x_2, \nabla f(x_2))$
How can we estimate $\nabla^2 f(x)$?
- In 1D:

$$\nabla^2 f(x) \approx \frac{\nabla f(x_2) - \nabla f(x_1)}{x_2 - x_1}$$

- In \mathbb{R}^n : let $y = \nabla f(x_2) - \nabla f(x_1)$, $\delta = x_2 - x_1$

Find “simplest” approximate Hessian matrix H or H^{-1} to fulfil

$$H \delta \stackrel{!}{=} y \quad \text{or} \quad \delta \stackrel{!}{=} H^{-1}y \quad (1)$$

(The first equation is called *secant equation*)

Basic example

- We've seen already two data points $(x_1, \nabla f(x_1))$ and $(x_2, \nabla f(x_2))$
How can we estimate $\nabla^2 f(x)$?

- In 1D:

$$\nabla^2 f(x) \approx \frac{\nabla f(x_2) - \nabla f(x_1)}{x_2 - x_1}$$

- In \mathbb{R}^n : let $y = \nabla f(x_2) - \nabla f(x_1)$, $\delta = x_2 - x_1$

Find “simplest” approximate Hessian matrix H or H^{-1} to fulfil

$$H \delta \stackrel{!}{=} y \quad \text{or} \quad \delta \stackrel{!}{=} H^{-1} y \quad (1)$$

(The first equation is called *secant equation*)

- Symmetric rank-1 solutions for H and H^{-1} :

$$H = \frac{yy^\top}{y^\top \delta} \quad \text{or} \quad H^{-1} = \frac{\delta \delta^\top}{\delta^\top y} \quad (2)$$

[Left: how to update H . Right: how to update directly H^{-1} .] Unconstrained Optimization Basics – Quasi-Newton methods – 28/41

BFGS

- Broyden-Fletcher-Goldfarb-Shanno (BFGS) method:

Input: initial $x \in \mathbb{R}^n$, functions $f(x), \nabla f(x)$, tolerance θ

Output: x

1: initialize $H^{-1} = \mathbf{I}_n$

2: **repeat**

3: compute $\delta = -H^{-1} \nabla f(x)$

4: perform a line search $\min_{\alpha} f(x + \alpha \delta)$

5: $\delta \leftarrow \alpha \delta$

6: $y \leftarrow \nabla f(x + \delta) - \nabla f(x)$

7: $x \leftarrow x + \delta$

8: update $H^{-1} \leftarrow \left(\mathbf{I} - \frac{y \delta^{\top}}{\delta^{\top} y} \right)^{\top} H^{-1} \left(\mathbf{I} - \frac{y \delta^{\top}}{\delta^{\top} y} \right) + \frac{\delta \delta^{\top}}{\delta^{\top} y}$

9: **until** $\|\delta\|_{\infty} < \theta$

- Notes:

– The blue term is the H^{-1} -update as on the previous slide

– The red term “deletes” previous H^{-1} -components. Check: $H^{-1} y = \delta$

– equivalent to the Sherman-Morrison formula: $H \leftarrow H - \frac{H \delta \delta^{\top} H^{\top}}{\delta^{\top} H \delta} + \frac{y y^{\top}}{y^{\top} \delta}$

L-BFGS

- In high dimensions, we do not want to explicitly store a dense H^{-1} . Instead we store vectors $\{v_i\}$ such that $H^{-1} = \sum_i v_i v_i^\top$
- **L-BFGS** (limited memory BFGS) limits the rank of the H^{-1} and thereby the used memory (number of vectors v_i)

L-BFGS

- In high dimensions, we do not want to explicitly store a dense H^{-1} . Instead we store vectors $\{v_i\}$ such that $H^{-1} = \sum_i v_i v_i^\top$
- **L-BFGS** (limited memory BFGS) limits the rank of the H^{-1} and thereby the used memory (number of vectors v_i)
- Some thoughts:
In principle, there are alternative ways to estimate H^{-1} from the data $\{(x_i, f(x_i), \nabla f(x_i))\}_{i=1}^k$, e.g. using Gaussian Process regression with derivative observations
 - not only the derivatives but also the value $f(x_i)$ should give information on $H(x)$ for non-quadratic functions
 - should one weight ‘local’ data stronger than ‘far away’? (GP covariance function)
 - related to model-based search (see Blackbox Optimization lecture)

(Nonlinear) Conjugate Gradient

Conjugate Gradient

- The “Conjugate Gradient Method” is a method for solving (large, or sparse) linear eqn. systems $Ax + b = 0$, without inverting or decomposing A . The steps will be “ A -orthogonal” (=conjugate). We mention its extension for optimizing nonlinear functions $f(x)$
 - As before we evaluated $g' = \nabla f(x_1)$ and $g = \nabla f(x_2)$ at two different points $x_1, x_2 \in \mathbb{R}^n$
 - Additional assumption: *exact line-search* step to x_2 :
 - $x_2 = x_1 + \alpha\delta_1$, $\alpha = \operatorname{argmin}_{\alpha} f(x_1 + \alpha\delta_1)$
 - iso-lines of $f(x)$ at x_2 are tangential to δ_1
- ⇒ The next search direction should be “orthogonal” to the previous one, but orthogonal w.r.t. the Hessian H , i.e., $d_2^\top H d_1 = 0$, which is called conjugate

Conjugate Gradient

Input: initial $x \in \mathbb{R}^n$, functions $f(x)$, $\nabla f(x)$, tolerance θ

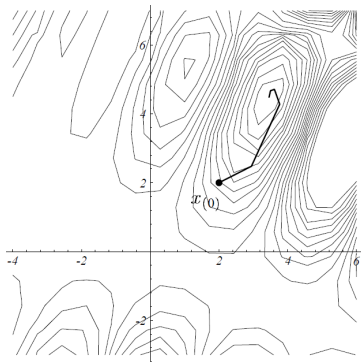
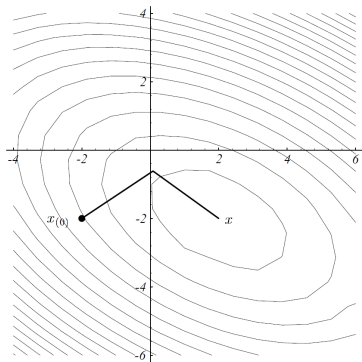
Output: x

- 1: initialize descent direction $d = g = -\nabla f(x)$
 - 2: **repeat**
 - 3: $\alpha \leftarrow \operatorname{argmin}_{\alpha} f(x + \alpha d)$ *// line search*
 - 4: $x \leftarrow x + \alpha d$
 - 5: $g' \leftarrow g$, $g = -\nabla f(x)$ *// store and compute grad*
 - 6: $\beta \leftarrow \max \left\{ \frac{g^{\top}(g-g')}{g^{\top}g'}, 0 \right\}$
 - 7: $d \leftarrow g + \beta d$ *// conjugate descent direction*
 - 8: **until** $|\Delta x| < \theta$
-

- **Notes:**

- $\beta > 0$: The new descent direction always adds a bit of the old direction!
- This *momentum* essentially provides 2nd order information
- The equation for β is by Polak-Ribière: On a quadratic function $f(x) = x^{\top}Ax + b^{\top}x$ this leads to **conjugate** search directions, $d^{\top}Ad = 0$.

Conjugate Gradient



- For quadratic functions CG converges in n iterations. But each iteration does *line search*

Rprop

Rprop

“Resilient Back Propagation” (outdated name from NN times...)

Input: initial $x \in \mathbb{R}^n$, function $f(x)$, $\nabla f(x)$, initial stepsize α , tolerance θ

Output: x

```
1: initialize  $x = x_0$ , all  $\alpha_i = \alpha$ , all  $g_i = 0$ 
2: repeat
3:    $g \leftarrow \nabla f(x)$ 
4:    $x' \leftarrow x$ 
5:   for  $i = 1 : n$  do
6:     if  $g_i g'_i > 0$  then // same direction as last time
7:        $\alpha_i \leftarrow 1.2\alpha_i$ 
8:        $x_i \leftarrow x_i - \alpha_i \text{sign}(g_i)$ 
9:        $g'_i \leftarrow g_i$ 
10:    else if  $g_i g'_i < 0$  then // change of direction
11:       $\alpha_i \leftarrow 0.5\alpha_i$ 
12:       $x_i \leftarrow x_i - \alpha_i \text{sign}(g_i)$ 
13:       $g'_i \leftarrow 0$  // force last case next time
14:    else
15:       $x_i \leftarrow x_i - \alpha_i \text{sign}(g_i)$ 
16:       $g'_i \leftarrow g_i$ 
17:    end if
18:    optionally: cap  $\alpha_i \in [\alpha_{\min} x_i, \alpha_{\max} x_i]$ 
19:  end for
20: until  $|x' - x| < \theta$  for 10 iterations in sequence
```

Rprop

- Rprop is a bit crazy:
 - stepsize adaptation in each dimension *separately*
 - it not only ignores $|\nabla f|$ but also its exact direction
 - step directions may differ up to $< 90^\circ$ from ∇f
 - Often works very robustly
 - Guarantees? See work by Ch. Igel
- If you like, have a look at:
Christian Igel, Marc Toussaint, W. Weishui (2005): Rprop using the natural gradient compared to Levenberg-Marquardt optimization. In Trends and Applications in Constructive Approximation. International Series of Numerical Mathematics, volume 151, 259-272.

Bound Constrained Optimization – postponed

- Bound constrained optimization:

$$\min_x f(x) \quad \text{s.t.} \quad l_i \leq x_i \leq u_i$$

simple box constraints on the variables (“simple constraints”)

- It might seem straight-forward to extend methods, perhaps by just clipping steps – but turns out not at all straight-forward
→ we’ll discuss it after constrained optimization

(e.g., Bertsekas: “Projected Newton methods for optimization problems with simple constraints”)

Appendix

Stopping Criteria

- Standard references (Boyd) define stopping criteria based on the “change” in $f(x)$, e.g. $|\Delta f(x)| < \theta$ or $|\nabla f(x)| < \theta$.
- Throughout I will define stopping criteria based on the change in x , e.g. $|\Delta x| < \theta$ repeatedly. In my experience with certain applications this is more meaningful, and invariant to the scaling of f . But this is application dependent.

Evaluating optimization costs

- Standard references (Boyd) assume line search is cheap and measure optimization costs as the number of iterations (counting 1 per line search).
- Throughout I will assume that every evaluation of $f(x)$ or $(f(x), \nabla f(x))$ or $(f(x), \nabla f(x), \frac{\partial}{\partial x} \phi(x))$ is approx. equally expensive—as is the case in certain applications.