

Numerical Methods for Nonconvex Optimization

- that are applicable when f is smooth: TODAY
- that are applicable when f is nonsmooth: Lec 13.

So, assume f is smooth (but not necessarily convex).
Two classes of methods

LINE SEARCH (Armijo cond + Wolfe cond.)

TRUST REGION: won't discuss. See Nocedal + Wright (NW).

In all cases, looking for local minimizers: global optimization
in nonconvex case is ~~HARD~~!

LINE SEARCH METHODS

1. Newton's method: if $H = \nabla^2 f(x)$ is not positive definite
then $d = -H^{-1} \nabla f(x)$ may not be a descent direction:

$$d^T \nabla f(x) = -d^T H d \text{ which could be } +ve$$

What to do?

(a) compute it anyway*, & use convex comb. of Newton + gradient
directions which can choose to be a descent direction.

(b) use Cholesky to factor H , & if it breaks down (which means
 $H \neq 0$), add a multiple of I to H & repeat

(c) use "modified Cholesky factorization" of Gill, Murray, Wright.
See NW. Not available as a standard code in Matlab.

(d) use "eig" to compute eigenvalues of H , & replace
negative or zero eigenvalues by small +ve numbers. This
results in a direction dominated by the "negative curvature"
directions. But it's several times the cost of Cholesky.

* a factorization that exploits symmetry but does
not assume that $H \succ 0$ is the LDL^T factorization, where

L has 1×1 and 2×2 blocks (Bunch-Pardell-Kaufman).

See the "help \\" info in Matlab.

CAN USE BACKTRACKING LINE SEARCH FOR ALL OF THESE.

2. Quasi-Newton + Conjugate Gradient Methods.

For these need a more sophisticated "ARMJO-WOLFE" line search.

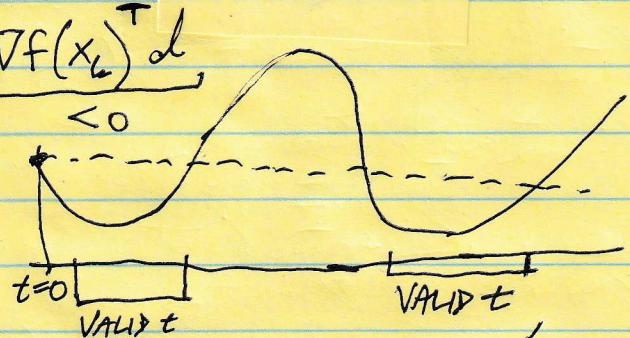
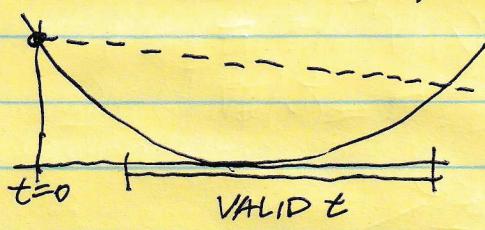
Recall the "Armijo condition". Let d = "search direction" Δx .

$$(A) \quad f(x_k + t d) \leq f(x_k) + c_1 t \nabla f(x_k)^T d$$

\leftarrow called α earlier, $\in (0, 1/2)$

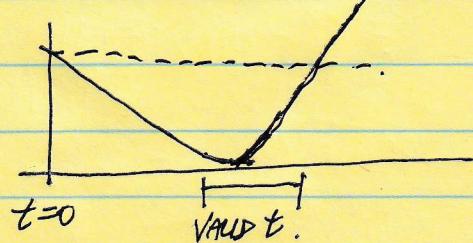
New condition: "Wolfe condition"

$$(W) \quad \nabla f(x_k + t d) \geq c_2 \underbrace{\nabla f(x_k)^T d}_{< 0}$$



(A) ensures t not too large

(W) ensures t not too small

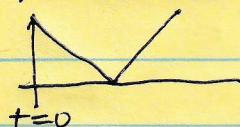


Require $0 < c_1 < c_2 < 1$.

Often, a "strong Wolfe" condition is recommended instead

$$(SW) \quad |\nabla f(x_k + t d)| \leq c_2 |\nabla f(x_k)^T d| \quad (\text{see NW})$$

But this can greatly reduce the range of acceptable steps, makes it potentially impossible if f is nonsmooth,



and makes writing a valid line search code MUCH MORE COMPLICATED: NW Ch. 3, section 6 is almost unreadable for this reason.

* Assume d is a DESCENT DIRECTION: $\nabla f(x_k)^T d < 0$.
(from x_k)

A simple Bracketing Armijo-Wolfe line search

let $t \leftarrow 1$, $\text{done} \leftarrow \text{false}$, $\alpha \leftarrow 0$, $\beta \leftarrow \infty$

UPPER BOUND

while not done

LOWER BOUND

```

 $x \leftarrow x_k + t d$  ( $d$  is "descent direction")
if  $f(x) > f(x_k) + c_1 t \nabla f(x_k)^T d$ 
     $\beta \leftarrow t$  : (A) condition violated,  $t$  too big
else if  $\nabla f(x)^T d < c_2 \nabla f(x_k)^T d$ 
     $\alpha \leftarrow t$  : (W) condition violated,  $t$  too small
else
     $\alpha \leftarrow t$ ,  $\beta \leftarrow t$ ,  $\text{done} \leftarrow \text{true}$ 
end.

```

(Set up next function evaluation)

```

if  $\beta < \alpha$ 
     $t \leftarrow (\alpha + \beta)/2$  (bisection)
else
     $t \leftarrow 2t$  (double)
end

```

code:

lineach_ww.m

in HANSO

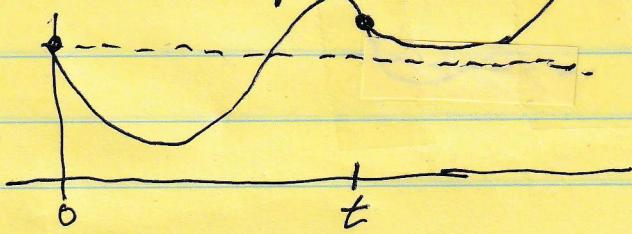
(see my
web page)
"weak
Wolfe" (N&W)

= "Carmijo"
"Wolfe".

Update chosen
so $[\alpha, \beta]$ always
contains a valid
A-W step t .

Essential that violation of 1st condition is checked first

e.g.



both conditions violated.

Update upper bound β , not lower bound α

This idea goes back to Powell (1976).

Thm If f is C^1 and bounded below on $\{x_k + t d : t \geq 0\}$,
then the line search terminates with a valid A-W step.

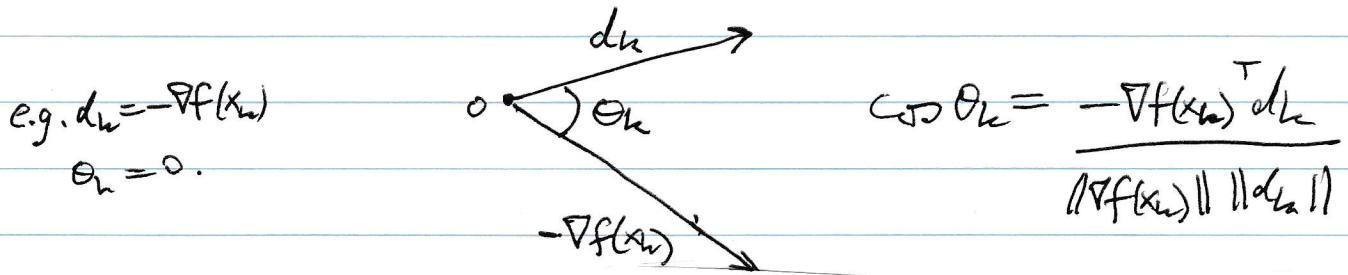
For an analysis under much weaker assumptions on f ,
e.g. f locally Lipschitz but not nec. differentiable, see LEWIS &
OVERTON 2013.

ZOUTENDIJK'S THEOREM

Assume f is bd below, $f \in C^1$ and ∇f is Lipschitz on $\{x : f(x) \leq f(x_0)\}$. (lip const. L)

Define a descent algorithm:

$$x_{k+1} = x_k + t_k d_k$$



where t_k satisfies the Armijo-Wolfe conditions.

Then as

$$\sum_{k=0}^{\infty} (\cos \theta_k)^2 \|\nabla f(x_k)\|^2 < \infty \quad (*)$$

so, if $\cos \theta_k \geq \tau > 0$ small k

$$(\theta_k \leq \psi < \frac{\pi}{2})$$

we have $\nabla f(x_k) \rightarrow 0$.

If write g_k for $\nabla f(x_k)$, f_k for $f(x_k)$.

We have from the Wolfe condition that

$$g_{k+1}^T d_k \geq c_2 g_k^T d_k$$

so $(g_{k+1} - g_k)^T d_k \geq (c_2 - 1) g_k^T d_k$

so $-t_k d_k \parallel d_k \geq (c_2 - 1) g_k^T d_k$

so $t_k \geq \frac{c_2 - 1}{L} g_k^T d_k / \|d_k\|^2$ (product of 2 neg. numbers)

Substitute this into the Armijo condition:

$$f_{k+1} \leq f_k + c_1 \underbrace{\frac{(c_2 - 1) g_k^T d_k}{L \|d_k\|^2}}_{\text{lowered onto } < 0} \cdot \frac{\|g_k\|^2}{\|g_k\|^2}$$

so

$$f_{k+1} \leq f_k - K (\cos \theta_k)^2 \cdot \frac{\|g_k\|^2}{\|g_k\|^2}$$

$\frac{c_1(1-c_2)}{L}$

sum over $j \leq k$

$$f_{k+1} \leq f_0 - K \sum_{j=0}^k (\cos \theta_j)^2 \|g_j\|^2$$

Since f is bd below, let $k \rightarrow \infty \Rightarrow (*)$.

Zoutendijk: this applies to any "Newton-like" method

$$d_k = H_k^{-1} \nabla f(x_k) \quad \text{as long as } H_k \text{ is}$$

UNIFORMLY POSITIVE DEF.

However, hard to prove this for QN methods —
and not true for CG methods.

Rate of convergence of Gradient Method: as in
convex case, show

— of Newton's Method: as before,
quadratic under regularity assumption, but
no guarantees can be made re tapers.

QUASI-NEWTON METHODS.

Motivation: Newton's method is $O(n^3)$ work.

Want to update an approx. to [FACTORIZATION OR INVERSE] of $\nabla^2 f(x)$ in $O(n^2)$ time.

How? Make better use of gradient info.

After line search, we have:

$$x_k \quad g_k \equiv \nabla f(x_k)$$

$$x_{k+1} \quad g_{k+1} = \nabla f(x_{k+1})$$

$$\text{Let } s_k = x_{k+1} - x_k = t_k d_k$$

$$y_k = g_{k+1} - g_k$$

From Fund. Thm. of Calc.,

$$y_k = \int_0^1 \nabla^2 f(x_k + \tau s_k) s_k d\tau$$

$$= \underbrace{\left[\int_0^1 \nabla^2 f(x_k + \tau s_k) d\tau \right]}_{G_k} s_k$$

G_k "average Hessian along s_k ".

so it seems reasonable that our new approx to $\nabla^2 f(x_{k+1})$, say B_{k+1} , should satisfy

$$B_{k+1} s_k = y_k$$

or, if we are approx $\nabla^2 f(x_{k+1})^{-1}$ by, say, C_k

THE
SECANT
EQUATION.

$$C_{k+1} y_k = s_k$$

Davidon-Fletcher-Powell
1959 1963.

11-7

Various choices known: PSB, DFP, BFGS

Broyden-Fletcher-Goldfarb-Shanno
1970

BFGS:

$$C_{k+1} = (I - \gamma_k s_k y_k^T) C_k (I - \gamma_k y_k s_k^T)^T + \gamma_k s_k s_k^T$$

$$\text{where } \gamma_k = \frac{1}{s_k^T y_k}$$

Check that: $C_{k+1} y_k =$

$$(I - \gamma_k s_k y_k^T) C_k (y_k - \underbrace{\gamma_k y_k s_k^T y_k}_0) + \gamma_k s_k s_k^T y_k \checkmark$$

How much work is needed to compute C_{k+1} ? ASK

⇒ How do we know $s_k^T y_k > 0$? DIRECTLY FROM THE WOLFE CONDITION (see next page).

THM (Powell, 1976)

Suppose 1. $f \in C^2$

2. $\Omega = \{x : f(x) \leq f(x_0)\}$ is convex

3. $\exists n, M \text{ s.t.}$

$$n \|z\|^2 \leq z^T \nabla f(x) z \leq M \|z\|^2$$

STRONG CONVEXITY

$\forall z \in \mathbb{R}^n, x \in \Omega$

Let $\{x_n\}$ generated by BFGS with Armijo-Wolfe line search

satisfies $x_n \rightarrow$ UNIQUE LOCAL MINIMIZER OF f . (existly).
(x^*)

Pf: beautiful, 2 pages + 30 minutes

See Nocedal + Wright.

HARD PART IS
SHOWING $\text{EIGS}(C_k)$
REMAIN BOUNDED,
AND $\lambda \geq \delta > 0$.

Clarify why $\mathbf{g}_k^T \mathbf{y}_w > 0$

11-7A.

$$\mathbf{g}_{k+1}^T \mathbf{d}_k \geq c_2 \underbrace{\mathbf{g}_k^T \mathbf{d}_k}_{< 0} \quad (\text{Wolfe})$$

$$\mathbf{y}_k^T \mathbf{d}_k = \mathbf{g}_{k+1}^T \mathbf{d}_k - \mathbf{g}_k^T \mathbf{d}_k \geq (c_2 - 1) \mathbf{g}_k^T \mathbf{d}_k > 0.$$

$$s_k = t_k d_k$$

$$\therefore \mathbf{s}_k^T \mathbf{y}_w = \mathbf{y}_w^T \mathbf{s}_k \geq (c_2 - 1) \mathbf{g}_k^T \mathbf{s}_k > 0.$$

Why update C_{k+1} , approx to $\nabla^2 f(\mathbf{x}_{k+1})^{-1}$, instead of B_{k+1} , approx to $\nabla^2 f(\mathbf{x}_{k+1})$. ASK.

Answer: Then $\mathbf{d}_{k+1} = C_{k+1} \nabla f(\mathbf{x}_{k+1})$ is only $O(n^2)$.

Alternative update Cholesky factor of $\nabla^2 f(\mathbf{x}_{k+1})$.

See Dennis + Schnabel's textbook. \rightarrow can solve for \mathbf{d}_{k+1} in $O(n^2)$ time.
Nice idea, used to be thought to be more stable, but not clear there is any advantage.

Superlinear Convergence (Dennis+Moré).

If we also assume $\nabla^2 f$ is Lipschitz, then

$x_n \rightarrow x^*$ superlinearly, i.e.

$$\boxed{\lim_{n \rightarrow \infty} \frac{\|x_{n+1} - x^*\|}{\|x_n - x^*\|} = 0}$$

Pf: See N&W.

NONCONVEX CASE: Nothing known in theory (if there are pathological counterexamples)
But in practice, always works! Like st. desc.

Methods that are $O(n)$

$\|\nabla f(x_n)\| \rightarrow 0$
CLUSTER POINTS
ARE STATIONARY
But unlike st. d., convergence is superlinear.

LIMITED MEMORY BFGS - far more efficient: $O(n)$
see Noc+Wri. But not

NONLINEAR CG let $d_0 = -\nabla f(x_0) = -g_0$. $\xrightarrow{\text{sup. conv.}}$

$$x_{k+1} = x_k + t_k d_k$$

$$d_{k+1} = -g_{k+1} + \beta_{k+1} d_k$$

$(-\nabla f(x_{k+1}))$

"AS IN
LINEAR"
CG

Hestenes-Rosens $\beta_{k+1}^{FR} = \frac{g_{k+1}^T g_{k+1}}{g_k^T g_k}$

Polak-Ribière

$$\beta_{k+1}^{PR} = \frac{g_{k+1}^T (g_{k+1} - g_k)}{g_k^T g_k}$$

In both cases reduce to linear "CG" when f is quadratic.

"Linear CG" \equiv solve $Ax = b$, $A \succ 0$

11-9

Theorem If $f(x) = x^T Ax - b^T x$ $A \succ 0$

and we use $t_h \leftarrow$ exact line search. (minimum of quadratic alongline)

then CG \equiv BFGS

and terminates in n steps. (More in presence of rounding
+ ill-conditioning of A),
OR FEWER! #DISTINCT EIG(A).

When n is large, want good approx in much
fewer than n steps.

Convergence is like $(1 - \sqrt{\frac{m}{M}})^k$ compared to $(1 - \frac{m}{M})^k$
in steepest descent.

Nonlinear Case (Nonquadratic)

To get convergence result for Nonlinear CG, need to
use a "strong Wolfe" line search to prove that FR
converges — but it's generally inferior to PR. (See Noeth Wini)

A variant: CGPRFR

$$\beta_{h+1} = \begin{cases} -\beta_{h+1}^{\text{FR}} & \text{if } \beta_{h+1}^{\text{PR}} < -\beta_{h+1}^{\text{FR}} \\ \beta_{h+1}^{\text{PR}} & \text{otherwise} \\ \beta_{h+1}^{\text{FR}} & \text{if } \beta_{h+1}^{\text{PR}} > \beta_{h+1}^{\text{FR}} \end{cases}$$

i.e. β_{h+1} is projection of β_{h+1}^{PR} onto $[-\beta_{h+1}^{\text{FR}}, \beta_{h+1}^{\text{FR}}]$

is a good compromise: as good or better than PR in practice,
same conv. theory as FR.

Recent work: variants that use only weak Wolfe.

Nonlinearly Constrained Optimization

A huge area.

Two big classes of algs.

SQP (sequential quadratic programming)

IP. (interior point)

See N+W.

Software SNOPT

IPOPT,

Margaret Wright will teach a course
that treats these in depth in Fall 2020.