

TURN ON RECORDING!

## NUMERICAL METHODS FOR NONCONVEX OPT.

TODAY: assume  $f$  is smooth

Sec 13: allow  $f$  to be nonsmooth

Looking for LOCAL minimizers.

LINE SEARCH METHODS - Today

TRUST REGION METHODS : N+L.

1. Newton's method.

If  $H = \nabla^2 f(x)$  is not pos def.,  
then  $d = -H^{-1} \nabla f(x)$  might not  
be a descent dir.

$$\nabla f(x)^T d = -\nabla f(x)^T H^{-1} \nabla f(x)$$

What to do?

(a) compute it anyway

could use a convex cont.

X  $d$  and  $\nabla f(x)$  which is a

(b) use Chuback +  $\overline{\text{descent dir?}}$

X if it breaks down, keep adding  
multiples of  $I$  until result is PD.

- (c) use "modified Cholesky fact" of Gill, Murray, Wright  
 - see N+W
- (d) use "eig" to get eigenvalues of  $H$ .
- (e) when chol breaks, you have  $R_i^T R_i = H_{ii}$   $H = \begin{bmatrix} H_{11} & \cdot \\ \cdot & I \end{bmatrix}$   
 Then solve  $\begin{bmatrix} H_{11} & 0 \\ 0 & I \end{bmatrix} d = -R_i^T f(x_i)$

## QUASI-NEWTON + NONLINEAR CONJUGATE GRADIENT.

Need an ARMJO-WOLFE line search.

Recall ARMJO cond:

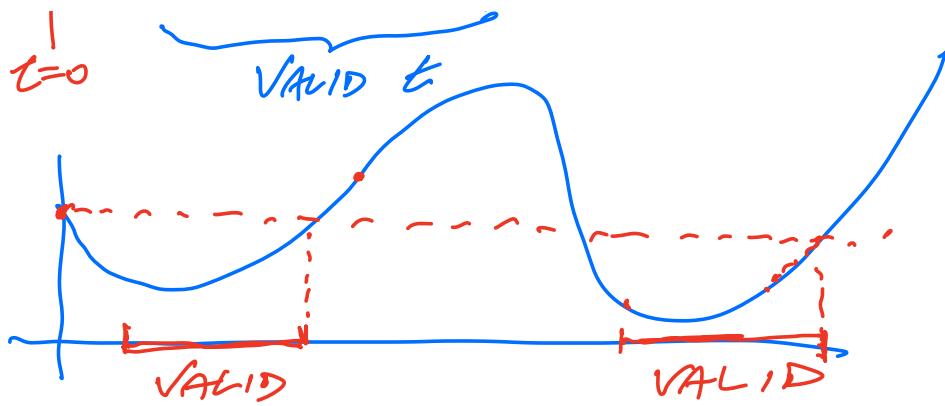
$$(A) f(x_i + t_i d_i) \leq f(x_i) + c_1 t_i \nabla f(x_i)^T d_i$$

and  $t_i$  before  
 $\in [0, 1/2]$

New: WOLFE cond:

$$(W) \quad \nabla f(x_i + t_i d_i)^T d_i \geq c_2 \nabla f(x_i)^T d_i$$

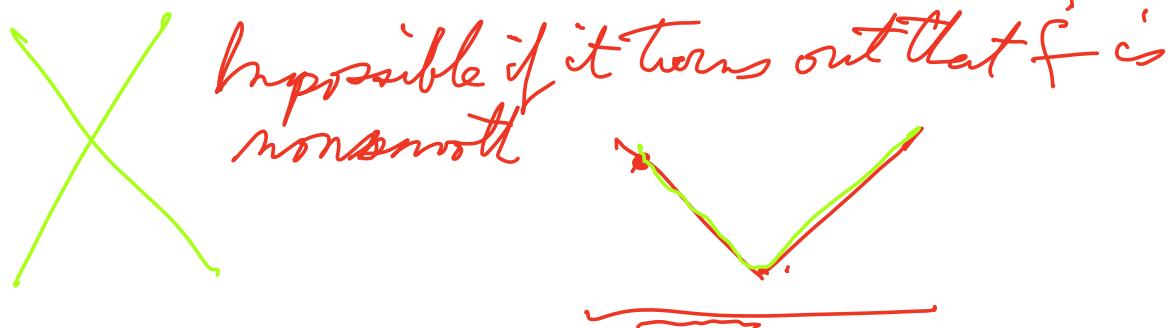
$c_1 < c_2 < 1$



(A) ensures that  $t$  not too large  
 (W) " " " " " small

Often a "strong Wolfe"

$$|\nabla f(x_k + t d_k)| \leq c_2 |\nabla f(x_k)|$$



A SIMPLE BRACKETING ARMIJO-WOLFE LINE SEARCH. Assume  $\nabla f(x)$  d/o

AWLS      set  $\epsilon \leftarrow 1$ , done = false,  $\alpha \leftarrow 0$ ,  $\beta \leftarrow \infty$

Property:  $[\alpha, \beta]$  always contains a valid A-W step  $t$

while not done

$$\bar{x} \leftarrow x_k + t d$$

$$\text{if } f(\bar{x}) > f(x_k) + c \cdot \epsilon \nabla f(x_k)^T d$$

```

t := t + d
 $\beta < t : (A) \text{ is violated}$   

 $\text{so } t \text{ is too large}$ 
else if  $\nabla f(x)^T d < c_2 \nabla f(x)^T d$ 
 $\alpha < t : (W) \text{ is violated}$   

 $\text{so } t \text{ is too small}$ 
else  $\alpha \leftarrow t, \beta \leftarrow t, \text{done} \leftarrow$   

end true  

% setup next f evaluation  

if  $\beta < \infty$   

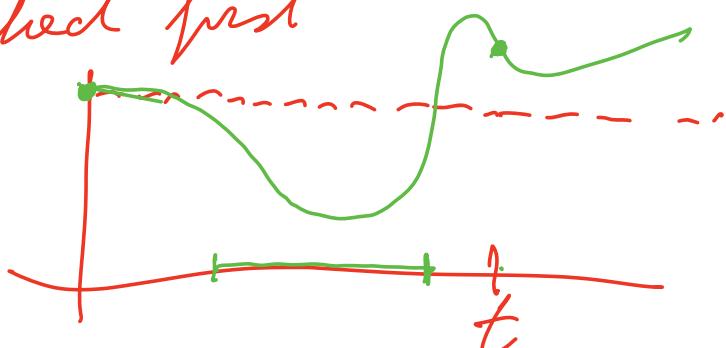
 $t \leftarrow \frac{\alpha+\beta}{2}$  (bisection)  

else  $t = 2t$  (double)  

end

```

Essential that violation of (A) is checked first



Then if  $f' \in C^1$  and bounded below,  
or  $\{x_n + td : t \geq 0\}$   
if  $x_n + td$  is  $\in \mathbb{R}$  for all  $t \geq 0$

then the A-W) generates w/n  
a valid A-W step.

For an analysis under weaker ass'ns on  
 $f$ , see Lewis & Overton 2013.

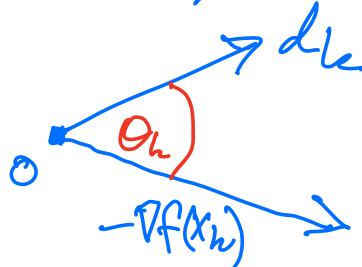
Code: linesch\_WW.m  
"weak Wolfe" = "Armijo-Wolfe"  
in HANSO.

### ZOUTENDIJK'S THEOREM

Assume  $f$  is bd below,  $C^1$ ,  $\nabla f(x)$   
is Lipschitz on  $\{x : f(x) \leq f(x_0)\}$

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

where  $t_k$  satisfies the (A-W) condition  
and



$$\cos \theta_k = - \frac{\nabla f(x_k)^T d_k}{\|\nabla f(x_k)\| \|d_k\|}$$

Then  $\sum_{k=0}^{\infty} (\cos \theta_k)^2 \frac{\|\nabla f(x_k)\|^2}{\|\nabla f(x_k)\| \|d_k\|} \leq \infty$

Pf  $\dots$   $\nabla f(x_k)$   $\nabla f(x_{k+1})$   $\dots$

LL Write  $g_h \ln Vf(x_h)$ , then to  $f(x_h)$   
We have from (W)

$$g_{h+1}^T d_h \geq c_2 g_h^T d_h$$

$$\Rightarrow (g_{h+1} - g_h)^T d_h \geq (c_2 - 1) g_h^T d_h$$

$$\Rightarrow L \|t_h d_h\| \|d_h\| \geq (c_2 - 1) g_h^T d_h$$

$$\Rightarrow t_h \geq \frac{c_2 - 1}{L} \frac{g_h^T d_h}{\|d_h\|^2}$$

subst. into (A):

$$f_{h+1} \leq f_h + t_h c_1 g_h^T d_h$$

$$f_{h+1} \leq f_h + \frac{c_2 - 1}{L} \frac{g_h^T d_h}{\|d_h\|^2} c_1 \frac{g_h^T d_h}{\|g_h\|^2}$$

$$f_{h+1} \leq f_h - K \frac{(c_1 \cos \theta_h)^2 \|g_h\|^2}{c_1 (1 - c_2)}$$

sum over  $j \leq k$

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

Since  $f_i$  is bounded

Let  $k \rightarrow \infty$

Apply to any Newton-like method

$$d_k = -\boxed{H_k^{-1}} \nabla f(x_k)$$

as long as  $H_k$  is UNIFORMLY POS DEF

## QUASI-NEWTON METHODS

("secant methods")

Motivation:

Newton's method is  $O(n^3)$  work  
Instead: update an approximation  
to FACTORIZATION or INVERSE of  
 $\nabla^2 f(x_k)$  is  $O(n^2)$  time.

How? Exploit gradient info.

After line search we have

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

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

Let  $s_k = x_{k+1} - x_k$

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

$$\begin{aligned}
 y_h &= \int_0^1 \nabla^2 f(x_h + \tau s_h) s_h \, d\tau \\
 &= \int_0^1 (\nabla^2 f(x_h + \tau s_h) \, d\tau) s_h
 \end{aligned}$$

$G_h$  "average Hessian along  $s_h$ "

So this motivates requiring our new approx to  $\nabla^2 f(x_{h+1})$ , say  $B_{h+1}$  to satisfy

$$B_{h+1} s_h = y_h$$

the SECANT (or QN) EQN.

or, if we approx  $\nabla^2 f(x_{h+1})$  by, say,  $H_{h+1}$

$$H_{h+1} y_h = s_h$$

Famous choices: PSB, DFP

$\uparrow$        $\uparrow$   
 Davidon 1959      Hestenes  
 Powell 1963

BFGS

Broyden, Fletcher, Goldfarb, Shanno  
1970

$$H_{k+1} = \boxed{(I - \gamma_k S_k y_k^T) H_k (I - \gamma_k y_k S_k^T)} + \boxed{\gamma_k S_k S_k^T}$$

where  $\gamma_k = \frac{1}{S_k^T y_k} > 0$

Can check  $H_k > 0 \Rightarrow H_{k+1} > 0$

Check

$$\boxed{H_{k+1} y_{k+1}} = \boxed{(I - \gamma_k S_k y_k^T) H_k} \boxed{y_k - \gamma_k y_k S_k^T y_k} + \boxed{\gamma_k S_k S_k^T y_k}$$

Computing  $H_{k+1}$  is  $O(n^2)$

$$\text{Setting } d_k = \boxed{H_{k+1} \nabla f(x_k)} \quad O(n^2)$$

If we approx  $B_{k+1} \approx \nabla f(x_k)$  instead,  
then we have to solve linear system.  
Or, could update Cholesky fact.

$$O(B_{k+1}) = O(n^2)$$

THM (Powell, 1976)

Suppose (i)  $f \in C^2$ , convex

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

(iii)  $\exists m, M$  s.t.

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

strongly convex

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

Then  $\{x_k\}$  from BFGS with A-W line search, satisfies

$x_k \rightarrow$  unique local  
minimizer of  $f$   $\star$

Here  $H_0$  is any pos. def. matx.

Why is  $s_k^T y_k \geq 0$ ?

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

$$\underbrace{y_k^T d_{k+1}}_{= g_{k+1}^T d_k - g_k^T d_k} \geq (c_2 - 1) \underbrace{g_k^T d_k}_{> 0} > 0$$

$$s_k = t_k d_k,$$

$$\text{or } s_k^T y_k = y_k^T s_k \geq (c_2 - 1) q_k^T s_k > 0.$$

Superlinear Convergence (Dennis+More)

If also assume  $D^2 f$  is Lipschitz

on  $\Omega$  then  $x_k \rightarrow x^*$  superlinearly  
i.e.

$$\lim_{k \rightarrow \infty} \frac{\|x_{k+1} - x^*\|}{\|x_k - x^*\|} = 0$$

## LIMITED MEMORY BFGS

L-BFGS.  
(N+L, see 7.2)

Use  $O(n)$  storage and  $O(n)$  work.

Let  $m \leq n$

At iter k, we have  $x_k, Df(x_k)$  and  
we have SAVED from previous iterates

$$\{s_i, y_i\}, i=k-m, \dots, k-1$$

Choose  $H_k^0$ , usually  $t_k I$ , + set

$$H_k = \underbrace{V_{k-1}^T \cdots V_{k-m+1}^T}_{\text{red}} V_{k-m}^T H^0 V_{k-m} V_{k-m+1}^T \cdots V_{k-1}^T + \gamma_{k-m} V_{k-1}^T \cdots V_{k-m+1}^T S_{k-m} S_{k-m}^T V_{k-m+1} \cdots V_{k-1} + \dots + \gamma_{k-1} S_{k-1} S_{k-1}^T$$

Usually use "scaling":  $H_k^0 = T_k I$

where  $T_k = \frac{S_{k-1}^T y_{k-1}}{y_{k-1}^T y_{k-1}}$

Can show L-BFGS has linear convergence to minimizer when  $f$  is smooth + ~~convex~~ convex

### OTHER $O(n)$ METHODS

### NONLINEAR CONJUGATE GRADIENT

$$d_0 = -g_0 = -\nabla f(x_0)$$

for  $k=1, 2, \dots$

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

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

- - - - -

FLETCHER-REEVES

$$\beta_{k+1}^{\text{FR}} =$$

$$= \frac{\mathbf{g}_{k+1}^T \mathbf{g}_{k+1}}{\mathbf{g}_k^T \mathbf{g}_k}$$

POLAK-RIBIERE

$$\beta_{k+1}^{\text{PR}} =$$

$$= \frac{\mathbf{g}_{k+1}^T (\mathbf{g}_{k+1} - \mathbf{g}_k)}{\mathbf{g}_k^T \mathbf{g}_k}$$

In both cases, reduce to

"linear CG" for  $\mathbf{A}\mathbf{x} = \mathbf{b}$   $\mathbf{A} \succ 0$

$$\text{i.e. } \min \frac{1}{2} \mathbf{x}^T \mathbf{A} \mathbf{x} - \mathbf{b}^T \mathbf{x}$$

But so does BFGS.

A variant:

$\beta_{k+1}^{\text{FRPR}}$  is projection of  
 $\beta_{k+1}^{\text{PR}}$  onto  $[-\beta_{k+1}^{\text{FR}}, \beta_{k+1}^{\text{FR}}]$

## CONSTRAINED NONCONVEX OPTIMIZATION

Two excellent codes

SNOPT "successive quadratic programming"

IPOPT "interior point"

$$A = U \Sigma V^T$$

$$A^T A = V \Sigma \cancel{U^T} \Sigma V^T = V \Sigma^2 V^T$$