

# Optimization

## Classical problems:

- find  $x$  st  $f(x)$  is minimized

- interesting cases

$$x \in \mathbb{R}^n \quad \left[ \begin{array}{l} f \in C^2 \quad (\text{cont. and 1, 2 deriv}) \\ \text{cont} \\ f \in C^1 \\ f \in C^0, \quad f \text{ convex} \\ f \in C^0 \end{array} \right.$$

Constrained:

$$x \in \{u \mid g(u) = 0\}$$

Discrete:

$$x \in \{0, 1\}^n$$

General methods.:Search:

construct a seq  $x_i$   
 $st \quad x_N \rightarrow$  right answer

Qns:

- How?
- When to stop?
- Converge how fast?

Closed form:

- uncommon, but sometimes useful
- Write conditions that identify a solution, then find it

→ Very valuable in itself!

Cases and terminology:Continuous optimization

- usually,  $f$  is at least continuous
- ~~usually~~,  $x \in \mathbb{R}^n$   
for the moment

if for  $x \in B^n(x^*)$

↑  
open ball centered at  $x^*$

we have

$$f(x^*) \leq f(x)$$

then  $x^*$  is a local minimum

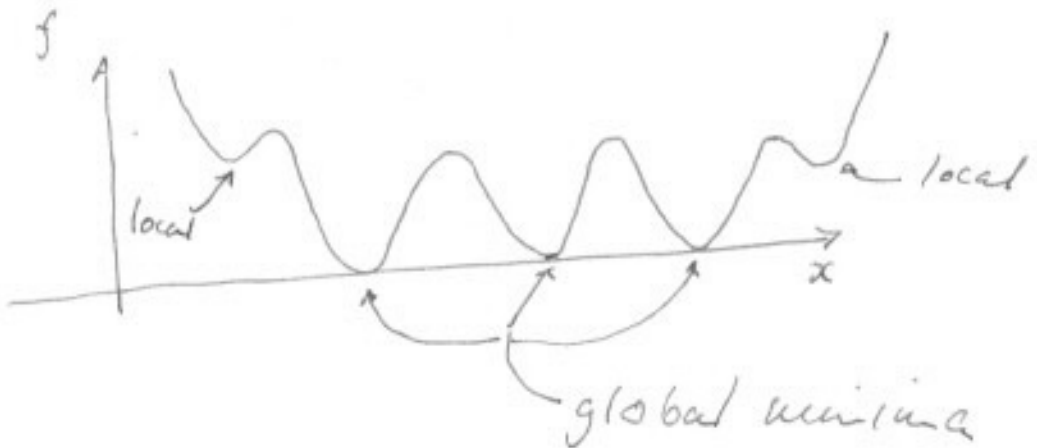
if we also have for  $x \in \mathbb{R}^n$

$$f(x^*) \leq f(x)$$

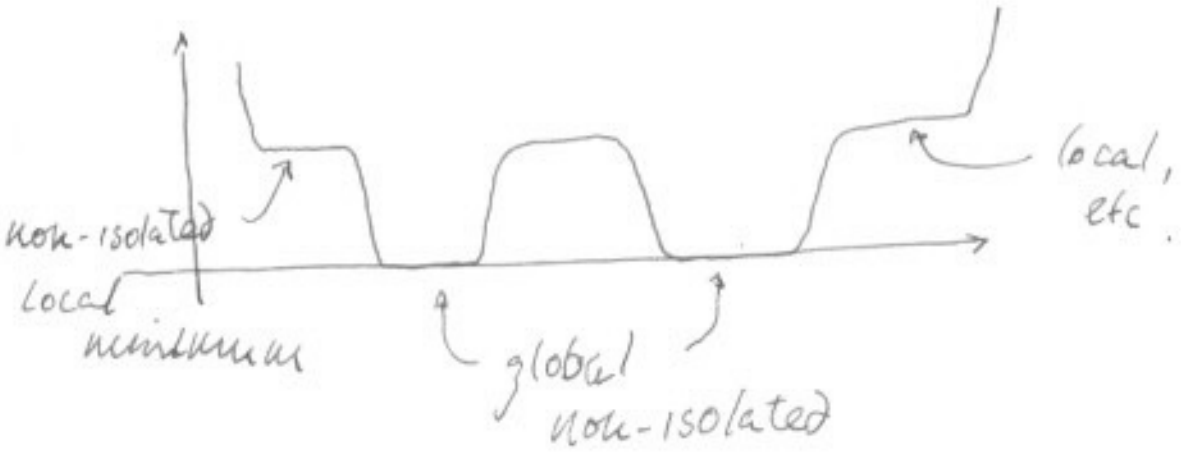
$x^*$  is a global minimum

Some examples suggest

- There can be many  $x^*$  that are a global minimum (But they all have the same  $f(x^*)$ !)



- Minima don't have to be isolated!

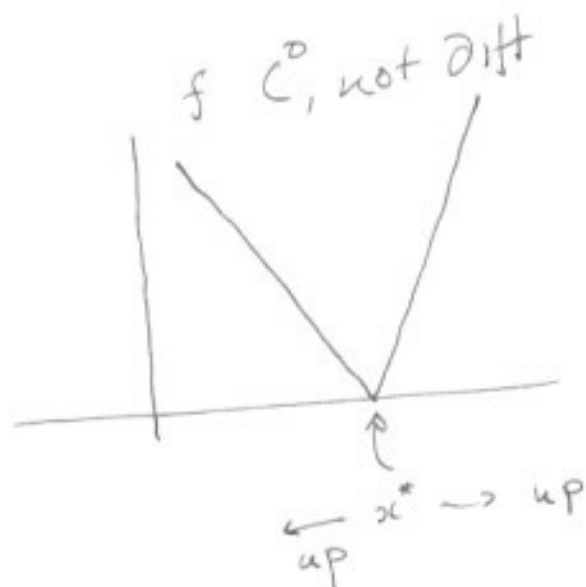
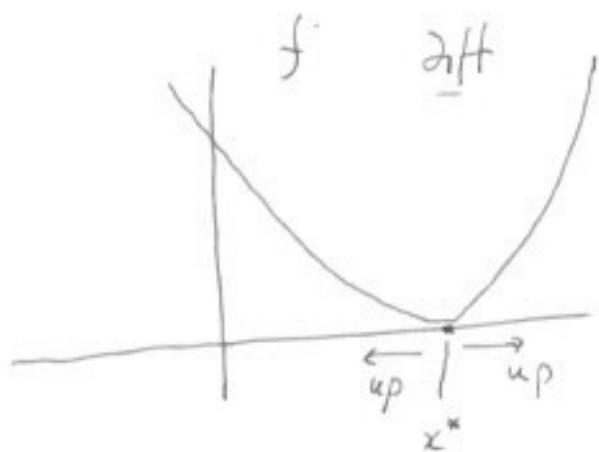


# Some closed form examples

① a - 4

General idea: if ~~at~~ any step in any direction takes you uphill you must be at a local min.

Most helpful when  $f$  is differentiable  
(but more to come)



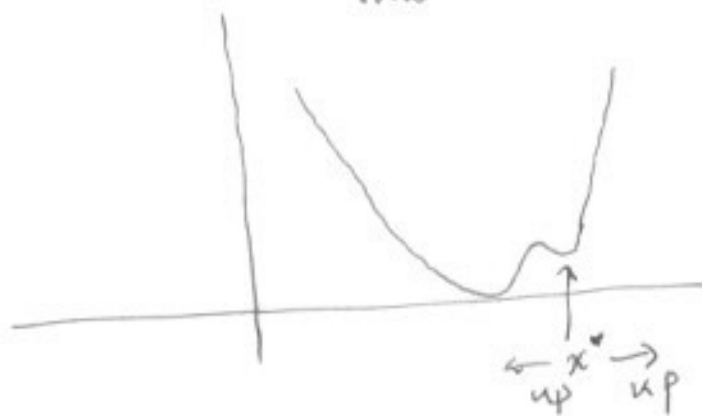
this means

$$f(x^* + \delta x) \geq f(x^*)$$

for any  $\delta x$  st.

$\|\delta x\|$  is sufficiently small.

↓  
this matters!



① a - 5.

Now assume  $f$  is continuously differentiable  
(i.e.  $f$  is continuous; derivative of  $f$  exists and  
is continuous)

$$f(x^* + \delta x) = f(x^*) + \nabla f^T \delta x + O(\delta x^2).$$

But at  $x^*$ , any direction is uphill

$$\text{So } \nabla f = 0$$

for small enough  
 $\|\delta x\|^2$

and this is useful.

Recipe 1:

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

$$\nabla f = Ax + b$$

→ Where and when is there a local minimum?

$$x \mid Ax + b = 0$$

↳ These are our suspects.

But these might not be  
local minima

① a-6

change coordinates:

$$u = x + w$$

where

$$Aw = b$$

does this  $w$  exist?  
what if it doesn't?

then

$$Au = 0$$

is our condition

Cases

A has full rank

→

$$u = 0$$

Q: ~~are~~ is  $\delta u^T A \delta u > 0$  for  
all small enough  $\delta u$ ?

i) If  $A \succ 0$ , yes.

positive definite;  
means  $w^T A w > 0$  for all  $w \neq 0$ .

ii) Not isolated local min!

A does not have full rank  
 $\rightarrow$  not isolated local min.

Remember this

$f(x) = \frac{x^T A x}{2} + b^T x + c$   
 has isolated local min at  
 $x$  st.  $Ax + b = 0$

if  $A \succ 0$   
 has isolated local max at  
 $x$  st.  $Ax + b = 0$

if  $-A \succ 0$



## Variational:

choose  $f$  such that

$$\int F(x, f(x)) dx \text{ is minimized}$$

How do we know we are at a local min?

- Any step in any direction makes cost get bigger
- Easy test if  $f$  is ~~differentiable~~  $C^1$

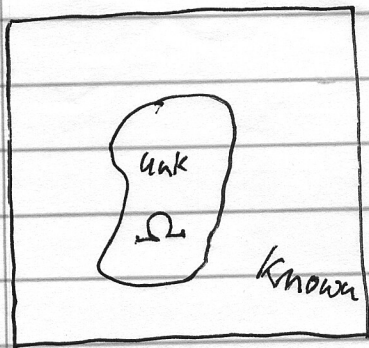
$$\nabla f = 0$$

- life is harder if  $f$  is  $C^0$  or worse

- local test might be difficult +

- eg.





we want to fill in  $\Omega$

reasonable criterion:

$$\min_{\Omega} \int \|\nabla f\|^2 dA$$

subject to:  $f = I$  on  $\partial\Omega$

i.e

- don't create derivatives unnecessarily
- agree w/ boundary.

• How could we solve this?

a)

- discretize, work with discretized derivative and integral
- we are now minimizing a function of a (big!) vector

Alternative:

• What properties does  $f$  have?

→ assume that  $\hat{f}$  is the soln.

→ now, for ANY test function  $\phi$ ,  
such that  $\phi = 0$  on  $\partial\Omega$ ,  
we have

$$\int_{\Omega} \|\nabla(f + \varepsilon\phi)\|^2 dA \geq \int_{\Omega} \|\nabla f\|^2 dA$$

for small enough  $\varepsilon > 0$

(i.e. if you make a small move in any  
direction, the value goes UP)

① - 6 - 3

Now, this means

$$\frac{d}{d\varepsilon} \int_{\Omega} \|\nabla f + \varepsilon \phi\|^2 dA = 0 \quad \text{for any } \phi$$

now this is

$$2 \int_{\Omega} \nabla f \cdot \nabla \phi dA = 0 \quad \left( \begin{array}{l} \text{doesn't seem} \\ \text{helpful} \end{array} \right)$$

But recall

$$\nabla \cdot [g \underline{v}] = \left( \nabla g \right) \cdot \underline{v} + g (\nabla \cdot \underline{v})$$

$$\text{i.e.} \quad \int_{\Omega} \nabla \cdot [\phi \nabla f] dA = \int_{\Omega} \nabla \phi \cdot \nabla f dA + \int_{\Omega} \phi (\nabla^2 f) dA$$

now

$$\int_{\Omega} \nabla \cdot [\phi \nabla f] dA = \int_{\partial \Omega} (\phi \nabla f) \cdot ds$$

(divergence theorem - remember!)

but  $\phi = 0$  on  $\partial \Omega$ .

$$\text{so } \int_{\Omega} \nabla \phi \cdot \nabla f dA = - \int_{\Omega} \phi (\nabla^2 f) dA = 0$$

But this is true for any  $\phi$ .

$$\text{so } \nabla^2 f = 0$$

and this offers other ways to solve.

① c

- Gives a nice criterion for variational case

$$\frac{d}{d\varepsilon} \left[ \int F(u, g(u) + \varepsilon \phi(u)) du \right] \Big|_{\varepsilon=0} = 0$$

Variational example



$$\min \int_0^a \sqrt{1 + \left(\frac{dg}{du}\right)^2} du$$

Subject to:

$$g(0) = 0$$
$$g(a) = b$$

i.e. for any test function  $\varphi$ ,  $\varphi(0) = 0$ ,  $\varphi(a) = 0$  ① d  
 we know  $g$  is right if

$$\frac{d}{d\varepsilon} \left[ \int \sqrt{1 + \left[ \frac{d}{dx} (g + \varepsilon \varphi) \right]^2} dx \right] \Big|_{\varepsilon=0} = 0$$

now write  $\frac{dg}{du} = g'$ ,  $\frac{d\varphi}{du} = \varphi'$

i.e.

Criterion is:

$$\int \frac{g' \varphi'}{\sqrt{1 + g'^2}} du = 0$$

Not promising; but use integration by parts

(i.e.  $\int u v' dx = uv \Big|_{x=0}^{x=a} - \int v u' dx$ )

① e

to get

$$\int \frac{d}{du} \left[ \frac{g'}{(1+g'^2)^{1/2}} \right] \cdot \varphi \, du = 0$$

for any  $\varphi$  ( $\varphi(0) = \varphi(a) = 0$ )

this gives  $\frac{d}{du} \left[ \frac{g'}{(1+g'^2)^{1/2}} \right] = 0$

which means

$$g'' \left[ \frac{1 - \frac{g'^2}{(1+g'^2)^{1/2}}}{(1+g'^2)} \right] = 0$$

and  $[\ ]$  is always positive so  $g'' = 0$

(The shortest distance between 2 points is a line!)

Equations obtained like this are

Euler Lagrange equations



14

Variational problems can be quite delicate

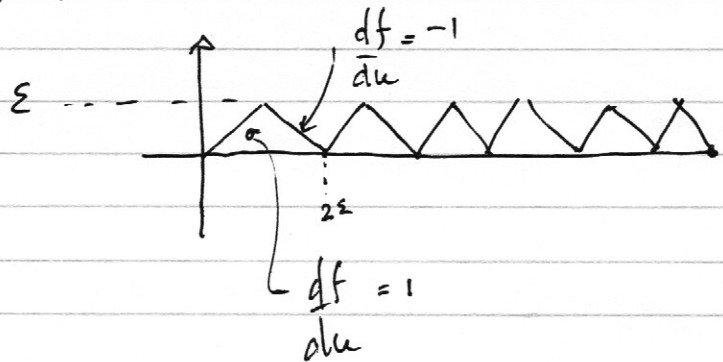
- a solution could not exist in reasonable function spaces.

example

$$\min_f \int_0^1 f(u)^2 + \left[ \left( \frac{df}{du} \right)^2 - 1 \right]^2 du = J[f]$$

$$\text{st. } f(0) = 0 ; f(1) = 0$$

Solu looks like



but as  $\epsilon$  gets smaller, we have

that  $J$  gets smaller, too

→ but there can't be a

limit

So no solution.

Failure of a solution to exist  
occurs in practical problems

19

• Simple issues:

$$\begin{array}{ccc} \operatorname{argmin}_u & u^2 & u \in (0, 1] \\ & & \uparrow \\ & & \text{open.} \end{array}$$

(This doesn't happen all that often,  
but is worth keeping in mind)

• Harder issues.

$$\operatorname{argmin}_f \int_0^1 \left[ \frac{1}{\sqrt{1+(f')^2}} - \frac{1}{\sqrt{2}} \right]^2 dx.$$

$$\text{s.t. } f(0) = 0$$

$$f(1) = 0$$

(h)

This sort of thing turns up in shape from shading problems rather often.

Notice I can get a min of the objective if  $f'^2 = 1$ .

→ So, if  $f \in C^\infty$  no solution (there can't be a  $C^p$  function st  $f'^2 = 1$ ,  $f(0) = 0$ ,  $f(1) = 0$ )

→ if  $f \in C^0$ , too many solutions!

^, ~, ~~~~~ etc.

~~Now assume that~~

①i

In practice, we usually turn variational problems into continuous optimization problems by writing

$$f = \sum_i a_i g_i$$

↑ basis functions

then solving for  $a_i$

BUT

- bad stuff can happen if original problem is poorly set
- The reasoning comes in useful later

Crucial, Take Home point:

You are at a minimum if every available step is uphill

① +

Now consider:

$$\min_x f(x) \quad f \in C^2$$

• A Descent Direction  $\underline{d}$  has the property that

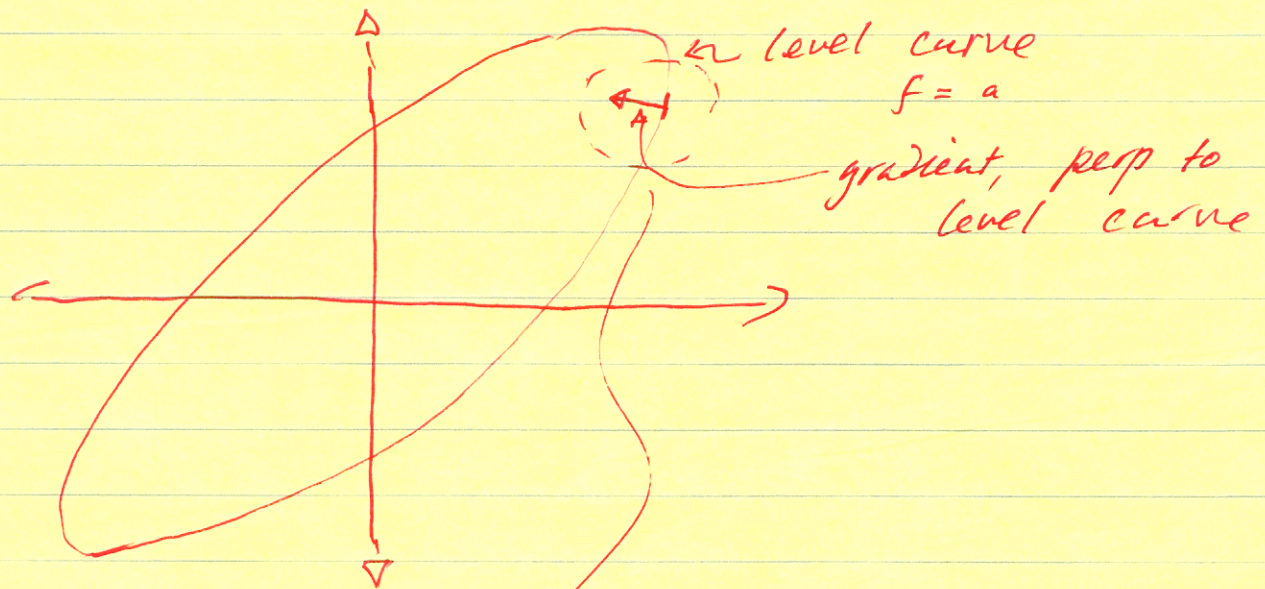
$$f(\underline{x}_0 + \varepsilon \underline{d}) < f(\underline{x}_0) \quad \text{for}$$

$$\varepsilon < \tau_d$$

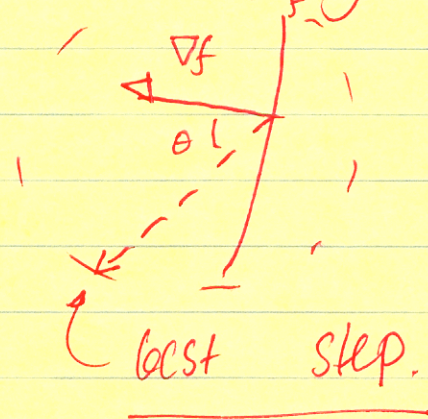
(3) a

Both gradient descent, coordinate descent are naughty:

- local model of function as quadratic form



→ best step is not along gradient



There are numerous descent directions

$$* \quad d_g = \frac{-\nabla f}{\|\nabla f\|}$$

this is gradient descent

issues:

- how to choose  $\epsilon$
- perhaps interval halving
- more sophisticated machinery later

\* write  $P_c$  for projection to some set of coordinate axes

- i.e.  $P_c$  zeros some elements

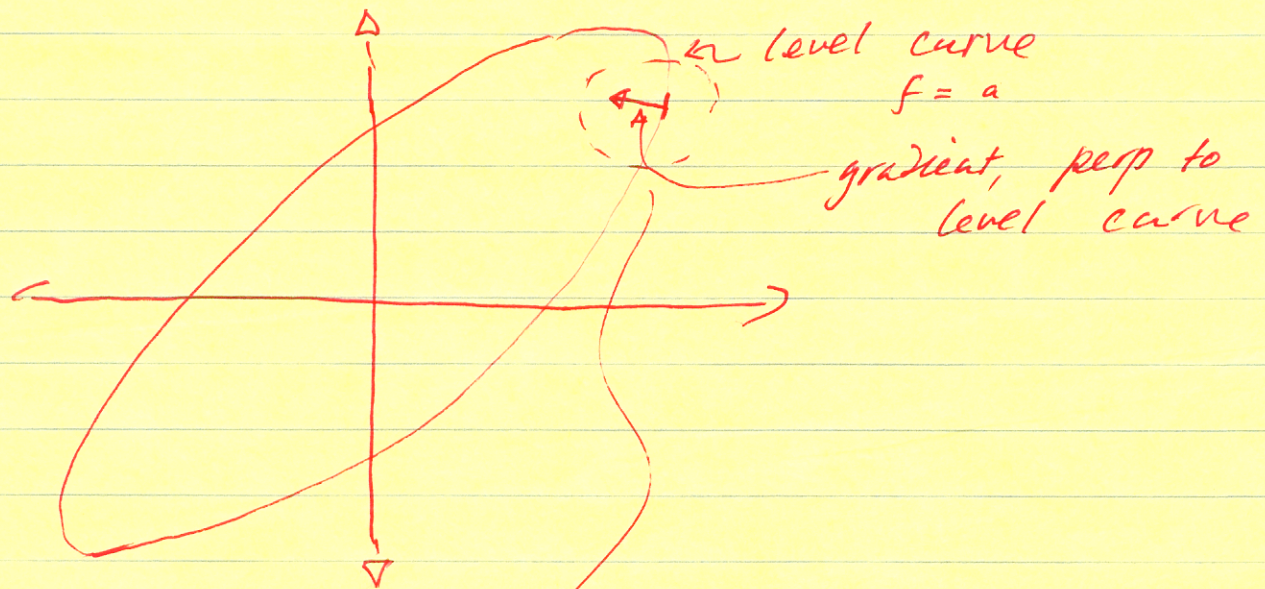
$$d_{cd} = -P_c d_g$$

this is coordinate descent

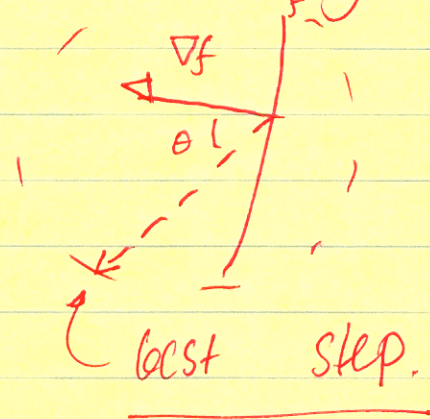
(3) a

Both gradient descent, coordinate descent are naughty:

- local model of function as quadratic form



→ best step is not  
along gradient





(3) 6

• if we take the best step along the gradient in this case, we don't go far  $\perp$  axis of symmetry



Now we zigzag slowly down the axis

Notice, best step can be at  $(-90^\circ, 90^\circ)$

to gradient

function is:

$$x' A x$$

, A positive definite

we are at  $u$ .

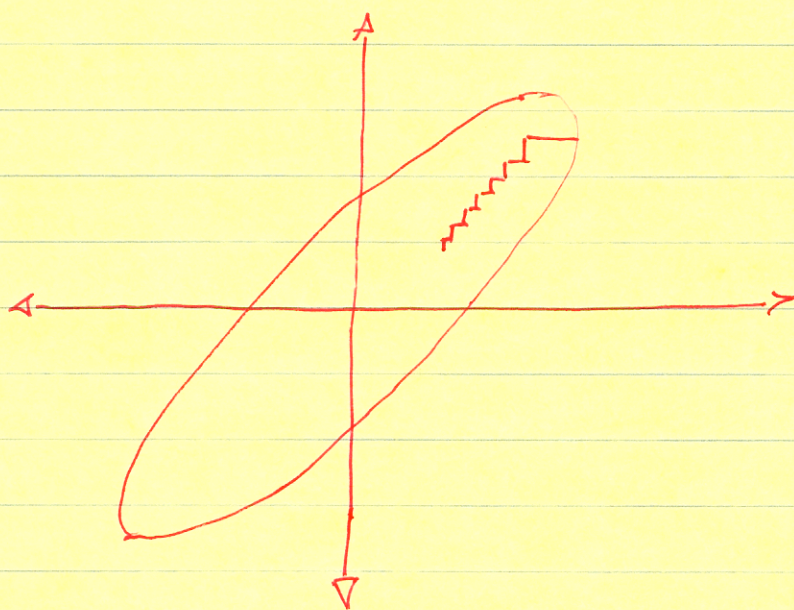
gradient is:  $2 A u$   
best step  $\lambda$  is  $-u$ .

(3)c

So cos angle is :

$$\frac{u'Au}{(u'u)^{1/2} (u'A'Au)^{1/2}}$$

Coordinate descent is also naughty



Newton's method:

$$f(x_0 + \delta) \approx f(x_0) + \nabla f' \delta + \frac{1}{2} \delta' H_f \delta + O(\delta^3)$$

- we could minimize the quadratic part as a function of  $\delta$ :

$$\min_{\delta} \quad \nabla f' \delta + \frac{1}{2} \delta' H_f \delta$$

i.e.  $\nabla f + H_f \delta = 0$

$$H_f \delta = -\nabla f$$

Q. Does Newton's method always give a descent direction?

$$f(x+d) = f(x) + \nabla f d + \frac{1}{2} d^T H_f d + O(d)^3$$

$$d = -H_f^{-1} \nabla f$$

$$f(x+d) - f(x) \approx -\frac{1}{2} \nabla f^T H_f^{-1} \nabla f$$

- so we're ok if  $H_f$  is positive definite

Q: is  $p$  a descent direction?

A: if  $p^T \nabla f < 0$

Notice we can obtain descent dir's from

$$p_k = -B_k^{-1} \nabla f_k$$

$\uparrow$   
iteration.

$$B_k^{-1} = Id \quad \text{gradient}$$

$$B_k^{-1} = P_c \quad \text{coord}$$

$$B_k^{-1} = H_f \quad \text{Newton.}$$

but

$B_k$  must be P.d for  $p$  to be a Descent Dirk.

Example: coordinate descent and EM

we have two parametric models

$$p_1(x|\theta_1) = e^{-g_1(x;\theta_1)}$$

and ~~the~~ ~~other~~ model

$$p_2(x|\theta_2) = e^{-g_2(x;\theta_2)}$$

and we observe  $x_i$  from a mixture

$$p(x|\theta) = \mu_1 p_1 + \mu_2 p_2$$

## Expectation - Maximization:

- Assume we have a mixture model

$$P(x|\theta) = \mu_1 P_1(x|\theta_1) + (1-\mu_1) P_2(x|\theta_2)$$

- for simplicity, I'll work with a mixture of exponentials

$$P_1(x|\theta_1) = e^{-g_1(x, \theta_1)}$$

$$P_2(x|\theta_2) = e^{-g_2(x, \theta_2)}$$

- I now have

$$x_1, \dots, x_n \sim P(x|\theta)$$

→ What is  $\theta = (\theta_1, \theta_2, \mu)$  ?

- Inference by maximum likelihood will be hard

because we must find

$$\arg \max_{\theta} \sum_i \log \left[ \mu_1 e^{-g_1(x_i, \theta_1)} + (1-\mu_1) e^{-g_2(x_i, \theta_2)} \right]$$

- This is difficult to work with, ~~an~~
  - The problem can be simplified if we know the mixture component from which each  $x_i$  comes.
 
$$S_i = \begin{cases} 1 & \text{if from 1} \\ 0 & \text{" " 2} \end{cases}$$
- ~~This gives~~  
 i.e.  $P(x_i | S_i, \theta)$  is straightforward

Write the Complete Data log-likelihood

$$\mathcal{J}(\theta) = \sum_i \log [P(x_i, S_i | \theta)]$$

$$= \sum_i \left[ \log P(x_i | S_i, \theta) + \log P(S_i | \theta) \right]$$

General algorithm:

$$h_c(\theta) = F(\delta, \theta)$$

we want to estimate these

↑ we don't know these.

assume we have an estimate  $\theta^{(n)}$  of  $\theta$   
then

$$Q(\theta; \theta^{(n)}) = E_{\delta | \theta^{(n)}, X} [F(\delta, \theta)]$$

function of  $\theta$  that  
depends on  $\theta^{(n)}$

E-Step:

Form  $Q(\theta; \theta^{(n)})$

M-Step

Find  $\theta^{(n+1)} = \underset{\theta}{\operatorname{argmax}} Q(\theta; \theta^{(n)})$



In our case

$$\begin{aligned} \mathcal{L}_c(\theta) &= \sum_i \left[ -\delta_i g_1(x_i, \theta_1) - (1-\delta_i) g_2(x_i, \theta_2) \right] \\ &\quad + \sum_i \left[ \delta_i \log \mu + (1-\delta_i) \log(1-\mu) \right] \end{aligned}$$

$E_{S_i | x, \hat{\theta}} \mathcal{L}_c(\theta)$  is straightforward, because

$\mathcal{L}_c(\theta)$  is linear in  $S_i$ .

→ Substitute  $P(S_i = 1 | x, \hat{\theta})$  for  $S_i$

$$P(S_i = 1 | x, \hat{\theta}) = \frac{P(x_i, S_i = 1 | \hat{\theta})}{P(x_i | \hat{\theta})}$$

$$= \frac{P(x_i | S_i = 1, \hat{\theta}) P(S_i = 1 | \hat{\theta})}{P(x_i | S_i = 1, \hat{\theta}) P(S_i = 1 | \hat{\theta}) +$$

$$P(x_i | S_i = 0, \hat{\theta}) P(S_i = 0 | \hat{\theta})}$$

In our case

$$P(\delta_i = 1 \mid x_i, \hat{\theta}^{(i)}) = \frac{e^{-g_1(x_i, \hat{\theta}_1^{(i)}) \hat{\mu}}}{e^{-g_1(x_i, \hat{\theta}_1^{(i)}) \hat{\mu}} + e^{-g_2(x_i, \hat{\theta}_2^{(i)}) (1-\hat{\mu})}}$$

Now: consider the following objective fn.

$$J(\theta, S) = J_c(\theta) + H(S)$$

↑  
Entropy

$$- \sum_i [\delta_i \log \delta_i + (1 - \delta_i) \log(1 - \delta_i)]$$

We will do coordinate descent

- fix  $\theta$ , min wrt  $\delta$
- fix  $\delta$ , min wrt  $\theta$

consider

$$\nabla_{\delta} \mathcal{F}(\theta, \delta) = 0$$

$$\frac{\partial \mathcal{F}}{\partial \delta_i} = - \left[ \log \delta_i - \log(1 - \delta_i) \right]$$

$$+ \left[ -g_1 + \log \hat{\mu} \right]$$

$$+ \left[ g_2 - \log(1 - \hat{\mu}) \right]$$

$$= 0$$

i.e.

$$\frac{s_i}{1-s_i} = \frac{e^{-g_1} \mu}{e^{-g_2} (1-\mu)}$$

and we substitute these in.

→ But this is what E step does.

$$\mu \text{ step} \equiv \nabla_{\theta} \hat{f}(\theta, \hat{s}_i) = 0$$

Q: Why not to Newton?

A: Not sure frankly,

H is big but sparse?

i.e.

$$\frac{\delta_i}{1 - \delta_i} = \frac{e^{-g_i} \mu}{e^{-g_2} (1 - \mu)}$$

hence: EM is coordinate ascent.

Q: why not so Newton's method?

A: not sure, frankly.

H is big, but sparse

Issues:

- What to do with a Zescent dirk?
- How to make H behave?
  - big
  - not P.D.
- How bad is gradient descent?

We have  $p_k$  and wish to choose a step length  $\alpha$ .

consider  $f(x_k + \alpha p_k)$   $\alpha > 0$

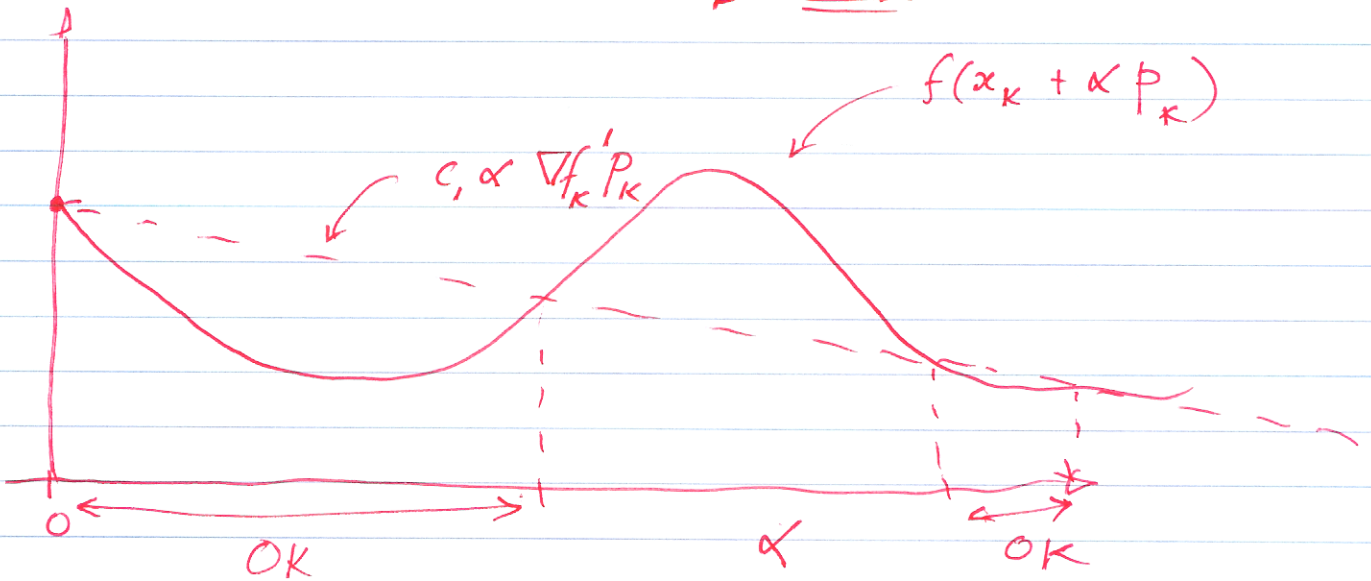
What  $\alpha$  are acceptable?

- ideally,  $\alpha$  is global minimizer
- sufficient decrease

$$f(x_k + \alpha p_k) \leq f(x_k) + c_1 \alpha \nabla f_k^T p_k$$

$$0 < c_1 < 1$$

for some constant  $\left[ \begin{array}{l} \text{Armijo condition} \\ \text{Wolfe} \end{array} \right]$   
 (typically  $10^{-4}$ )



Sufficient decrease is not enough

→ very small  $\alpha$  are OK.

$$\nabla f(x_k + \alpha p_k)^T p_k \geq c_2 \nabla f_k^T p_k$$

$$c_1 \leq c_2 \leq 1$$

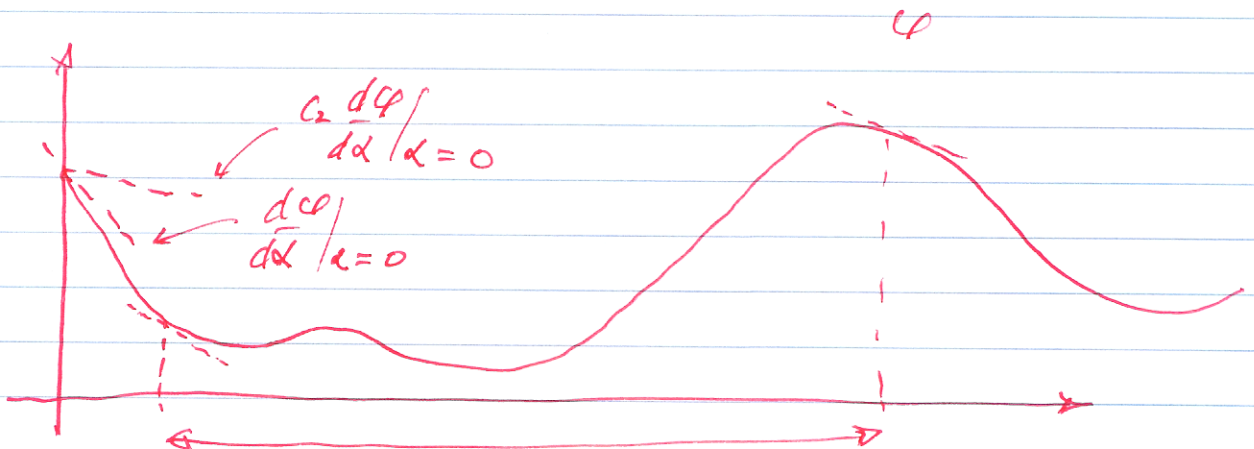
notice:

write  $\phi(\alpha) = f(x_k + \alpha p_k)$

then  $\frac{d\phi}{d\alpha} = \nabla f(x_k + \alpha p_k)^T p_k$

so condition is:

$$\frac{d\phi}{d\alpha} \geq c_2 \nabla f_k^T p_k = c_2 \left. \frac{d\phi}{d\alpha} \right|_{\alpha=0}$$



Notice sign of slope!

$c_2$  is usually 0.4 (Newton)  
0.1 (conj. grad.)

Wolfe conditions

$$f(x_k + \alpha_k p_k) \leq f(x_k) + c_1 \alpha_k \nabla f_k^T p_k$$

$$\nabla f(x_k + \alpha_k p_k)^T p_k \geq c_2 \nabla f_k^T p_k$$

Notice: for  $f$  continuously diff,  
 $f$  bounded below along  $x_k + \alpha p_k$ ,  $\alpha > 0$   
there exist intervals satisfying  
these conds.

Alg: for  $\tilde{\alpha} > 0$ ,  $\rho \in (0, 1)$   
 $\alpha = \tilde{\alpha}$   
repeat until [sufficient descent]  
 $\alpha = \rho \alpha$   
end

OK for Newton; not as good for others.



(12a)

Now we are generating a sequence  $\{x_i\}$  by finding  $P_k, \alpha_k$  and accepting

$$x_{k+1} = x_k + \alpha_k P_k$$

Q: How does this seq behave?

A:

Q: To what does it converge?

Q: How fast?

Some answers by defining

$$\cos \theta_k = \frac{-\nabla f_k^T P_k}{\|\nabla f_k\| \|P_k\|}$$

Thm: (Zoutendijk)

Consider iteration of given form,  $\alpha_k$  satisfying Wolfe cond,  $f$  bounded below continuously diff in an open set  $\Lambda$  containing  $L = \{x : f(x) \leq f(x_0)\}$ . Assume  $\nabla f$  is Lipschitz. Then

$$\sum_k \cos^2 \theta_k \|\nabla f_k\|^2 < \infty$$

(12c)

## Rates of Convergence:

- We have, for  $f \in C^2$ , exact line search (i.e. best  $\alpha_k$ ),  $x^*$  the min, Gradient descent behaves like

$$f(x_{k+1}) - f(x^*) \leq \rho^2 [f(x_k) - f(x^*)]$$

↑  
related to Hessian

- For Newton, if  $x_0$  sufficiently close to  $x^*$

$$\|x_{k+1} - x^*\| \leq L \|x_k - x^*\|^2$$

↑  
related to  
Hessian

Newton's method with Hessian modification

Problem:  $H$  may not be P.D.,

so

$H p_k = -\nabla f$  may not give

a descent direction

Strategy: modify  $H$  to be P.D.

$$B_k = H_f + E_k$$

↑  
chosen to make  $B_k$  PD

This will converge globally if

$\{H_f(x_k)\}$  is bounded

$\Rightarrow$

$$\|B_k\| \|B_k^{-1}\| \leq C > 0$$

(1A)

Generally would like  $K_k$  small  
(So as to preserve Hessian info)

1: Add a multiple of Identity:

choose  $\beta > 0$   
if  $\min_i h_{ii} > 0$

$$\tilde{\tau}_0 = 0$$

else

$$\tilde{\tau}_0 = -\min(h_{ii}) + \beta;$$

end.

for  $k = \dots$

attempt cholesky factorization of  $H + \tau_k I$   
if OK return factor

else

$$\tau_{k+1} = \max(2\tau_k, \beta)$$

end

end.

we are searching for  $\tau I$  to  
make  $H$  p.d.

Cholesky:

$$A = L L^T$$

lower triang

- works if A PD, otherwise

get a  $\sqrt{0}$ ,  $\sqrt{-ve}$

Modified Cholesky

$$A = L D L^T$$

where

L is lower triang, 1's on diag

D is diag, +ve diag

Note: if A pd, then D elements are +ve

Cholesky:

for  $j = 1 \dots n$

$$c_{jj} = a_{jj} - \sum_{s=1}^{j-1} d_s l_{js}^2$$

$$d_j = c_{jj}$$

for  $i = j+1 \dots n$

$$c_{ij} = a_{ij} - \sum_{s=1}^{j-1} d_s l_{is} l_{js}$$

$$l_{ij} = c_{ij} / d_j$$

end

end

Now:  $d_{jj}$  all positive if  $A$  PD.

Modify alg so that

$$d_{jj} = \max \left( |c_{jj}|, \left( \frac{\theta_j}{\beta} \right)^2, \delta \right)$$

$$\theta_j = \max_{j < i \leq n} |c_{ij}|$$

and this gives a "factorization"  
where

$$d_j \geq \delta$$
$$|m_{ij} = l_{ij} \sqrt{d_j}| \leq \beta$$

Desirable for  
error control.

### Improvements

- Permute rows and columns to reduce the size of the modification.
- This will give guaranteed bounds  $\Rightarrow$  global convergence.



Step length selection:

$$\varphi(\alpha) = f(x_0 + \alpha p_k)$$

Sufficient decrease is then,

$$\varphi(\alpha) \leq \varphi(0) + c_1 \alpha \varphi'(0).$$

• guess  $\alpha_0$

→ OK ; stop

→ Not OK ; there is an OK step in interval.

• we know  $\varphi(0)$ ,  $\varphi(\alpha_0)$ ,  $\varphi'(0)$

• build quadratic interpolate

$$\begin{aligned} \varphi(\alpha) &= \left( \frac{\varphi(\alpha_0) - \varphi(0) - \alpha_0 \varphi'(0)}{\alpha_0^2} \right) \alpha^2 \\ &\quad + \varphi'(0) \alpha \\ &\quad + \varphi(0) \end{aligned}$$

• minimize in  $\alpha$  to get  $\alpha_1$

→  $\alpha_1$  OK ; stop

→ else construct cubic

interpolate of  $\varphi(0)$   $\varphi'(0)$   $\varphi(\alpha_0)$

$\varphi(\alpha_1)$

minimize ;  $\alpha_2$

→  $\alpha_2$  OK stop

→ else cubic with  $\varphi(0)$ ,  $\varphi'(0)$ ,

two most recent  $\alpha$

It can be shown that if  $x_k \rightarrow x^*$  superlinearly, then the ratio in this expression converges to 1. If we adjust the choice (3.60) by setting

$$\alpha_0 \leftarrow \min(1, 1.01\alpha_0),$$

we find that the unit step length  $\alpha_0 = 1$  will eventually always be tried and accepted, and the superlinear convergence properties of Newton and quasi-Newton methods will be observed.

### A LINE SEARCH ALGORITHM FOR THE WOLFE CONDITIONS

The Wolfe (or strong Wolfe) conditions are among the most widely applicable and useful termination conditions. We now describe in some detail a one-dimensional search procedure that is guaranteed to find a step length satisfying the *strong* Wolfe conditions (3.7) for any parameters  $c_1$  and  $c_2$  satisfying  $0 < c_1 < c_2 < 1$ . As before, we assume that  $p$  is a descent direction and that  $f$  is bounded below along the direction  $p$ .

The algorithm has two stages. This first stage begins with a trial estimate  $\alpha_1$ , and keeps increasing it until it finds either an acceptable step length or an interval that brackets the desired step lengths. In the latter case, the second stage is invoked by calling a function called **zoom** (Algorithm 3.6, below), which successively decreases the size of the interval until an acceptable step length is identified.

A formal specification of the line search algorithm follows. We refer to (3.7a) as the *sufficient decrease condition* and to (3.7b) as the *curvature condition*. The parameter  $\alpha_{\max}$  is a user-supplied bound on the maximum step length allowed. The line search algorithm terminates with  $\alpha_*$  set to a step length that satisfies the strong Wolfe conditions.

#### Algorithm 3.5 (Line Search Algorithm).

Set  $\alpha_0 \leftarrow 0$ , choose  $\alpha_{\max} > 0$  and  $\alpha_1 \in (0, \alpha_{\max})$ ;

$i \leftarrow 1$ ;

**repeat**

Evaluate  $\phi(\alpha_i)$ ;

**if**  $\phi(\alpha_i) > \phi(0) + c_1\alpha_i\phi'(0)$  or  $[\phi(\alpha_i) \geq \phi(\alpha_{i-1})$  and  $i > 1]$

$\alpha_* \leftarrow \mathbf{zoom}(\alpha_{i-1}, \alpha_i)$  and **stop**;

Evaluate  $\phi'(\alpha_i)$ ;

**if**  $|\phi'(\alpha_i)| \leq -c_2\phi'(0)$

**set**  $\alpha_* \leftarrow \alpha_i$  and **stop**;

**if**  $\phi'(\alpha_i) \geq 0$

**set**  $\alpha_* \leftarrow \mathbf{zoom}(\alpha_i, \alpha_{i-1})$  and **stop**;

Choose  $\alpha_{i+1} \in (\alpha_i, \alpha_{\max})$ ;

$i \leftarrow i + 1$ ;

**end (repeat)**

Note that the se  
the order of the argu  
the knowledge that th  
conditions if one of th

(i)  $\alpha_i$  violates the s

(ii)  $\phi(\alpha_i) \geq \phi(\alpha_{i-1})$

(iii)  $\phi'(\alpha_i) \geq 0$ .

The last step of the al  
implement this step v  
we can simply set  $\alpha_{i+1}$   
important that the suc  
a finite number of iter

We now specify  
its input arguments is

(a) the interval bound  
conditions;

(b)  $\alpha_{lo}$  is, among all  
condition, the on

(c)  $\alpha_{hi}$  is chosen so th

Each iteration of **zoom**  
of these endpoints by c

#### Algorithm 3.6 (zoom)

**repeat**

Interpolate (u

a trial st

Evaluate  $\phi(\alpha_j)$ ;

**if**  $\phi(\alpha_j) > \phi(0)$

$\alpha_{hi} \leftarrow \alpha_j$

**else**

Evaluate

**if**  $|\phi'(\alpha_j)| \leq -c_2\phi'(0)$

**set**  $\alpha_* \leftarrow \alpha_j$

**if**  $\phi'(\alpha_j) \geq 0$

$\alpha_{lo} \leftarrow \alpha_j$

$\alpha_{hi} \leftarrow \alpha_j$

**end (repeat)**

Note that the sequence of trial step lengths  $\{\alpha_i\}$  is monotonically increasing, but that the order of the arguments supplied to the **zoom** function may vary. The procedure uses the knowledge that the interval  $(\alpha_{i-1}, \alpha_i)$  contains step lengths satisfying the strong Wolfe conditions if one of the following three conditions is satisfied:

- (i)  $\alpha_i$  violates the sufficient decrease condition;
- (ii)  $\phi(\alpha_i) \geq \phi(\alpha_{i-1})$ ;
- (iii)  $\phi'(\alpha_i) \geq 0$ .

The last step of the algorithm performs extrapolation to find the next trial value  $\alpha_{i+1}$ . To implement this step we can use approaches like the interpolation procedures above, or we can simply set  $\alpha_{i+1}$  to some constant multiple of  $\alpha_i$ . Whichever strategy we use, it is important that the successive steps increase quickly enough to reach the upper limit  $\alpha_{\max}$  in a finite number of iterations.

We now specify the function **zoom**, which requires a little explanation. The order of its input arguments is such that each call has the form **zoom** $(\alpha_{lo}, \alpha_{hi})$ , where

- (a) the interval bounded by  $\alpha_{lo}$  and  $\alpha_{hi}$  contains step lengths that satisfy the strong Wolfe conditions;
- (b)  $\alpha_{lo}$  is, among all step lengths generated so far and satisfying the sufficient decrease condition, the one giving the smallest function value; and
- (c)  $\alpha_{hi}$  is chosen so that  $\phi'(\alpha_{lo})(\alpha_{hi} - \alpha_{lo}) < 0$ .

Each iteration of **zoom** generates an iterate  $\alpha_j$  between  $\alpha_{lo}$  and  $\alpha_{hi}$ , and then replaces one of these endpoints by  $\alpha_j$  in such a way that the properties (a), (b), and (c) continue to hold.

**Algorithm 3.6** (zoom).

```

repeat
  Interpolate (using quadratic, cubic, or bisection) to find
    a trial step length  $\alpha_j$  between  $\alpha_{lo}$  and  $\alpha_{hi}$ ;
  Evaluate  $\phi(\alpha_j)$ ;
  if  $\phi(\alpha_j) > \phi(0) + c_1\alpha_j\phi'(0)$  or  $\phi(\alpha_j) \geq \phi(\alpha_{lo})$ 
     $\alpha_{hi} \leftarrow \alpha_j$ ;
  else
    Evaluate  $\phi'(\alpha_j)$ ;
    if  $|\phi'(\alpha_j)| \leq -c_2\phi'(0)$ 
      Set  $\alpha_* \leftarrow \alpha_j$  and stop;
    if  $\phi'(\alpha_j)(\alpha_{hi} - \alpha_{lo}) \geq 0$ 
       $\alpha_{hi} \leftarrow \alpha_{lo}$ ;
     $\alpha_{lo} \leftarrow \alpha_j$ ;
end (repeat)

```

If the new estimate  $\alpha_j$  happens to satisfy the strong Wolfe conditions, then **zoom** has served its purpose of identifying such a point, so it terminates with  $\alpha_* = \alpha_j$ . Otherwise, if  $\alpha_j$  satisfies the sufficient decrease condition and has a lower function value than  $x_{l_0}$ , then we set  $\alpha_{l_0} \leftarrow \alpha_j$  to maintain condition (b). If this setting results in a violation of condition (c), we remedy the situation by setting  $\alpha_{hi}$  to the old value of  $\alpha_{l_0}$ . Readers should sketch some graphs to see for themselves how **zoom** works!

As mentioned earlier, the interpolation step that determines  $\alpha_j$  should be safeguarded to ensure that the new step length is not too close to the endpoints of the interval. Practical line search algorithms also make use of the properties of the interpolating polynomials to make educated guesses of where the next step length should lie; see [39, 216]. A problem that can arise is that as the optimization algorithm approaches the solution, two consecutive function values  $f(x_k)$  and  $f(x_{k-1})$  may be indistinguishable in finite-precision arithmetic. Therefore, the line search must include a stopping test if it cannot attain a lower function value after a certain number (typically, ten) of trial step lengths. Some procedures also stop if the relative change in  $x$  is close to machine precision, or to some user-specified threshold.

A line search algorithm that incorporates all these features is difficult to code. We advocate the use of one of the several good software implementations available in the public domain. See Dennis and Schnabel [92], Lemaréchal [189], Fletcher [101], Moré and Thuente [216] (in particular), and Hager and Zhang [161].

One may ask how much more expensive it is to require the strong Wolfe conditions instead of the regular Wolfe conditions. Our experience suggests that for a “loose” line search (with parameters such as  $c_1 = 10^{-4}$  and  $c_2 = 0.9$ ), both strategies require a similar amount of work. The strong Wolfe conditions have the advantage that by decreasing  $c_2$  we can directly control the quality of the search, by forcing the accepted value of  $\alpha$  to lie closer to a local minimum. This feature is important in steepest descent or nonlinear conjugate gradient methods, and therefore a step selection routine that enforces the strong Wolfe conditions has wide applicability.

## NOTES AND REFERENCES

For an extensive discussion of line search termination conditions see Ortega and Rheinboldt [230]. Akaike [2] presents a probabilistic analysis of the steepest descent method with exact line searches on quadratic functions. He shows that when  $n > 2$ , the worst-case bound (3.29) can be expected to hold for most starting points. The case  $n = 2$  can be studied in closed form; see Bazaraa, Sherali, and Shetty [14]. Theorem 3.6 is due to Dennis and Moré.

Some line search methods (see Goldfarb [132] and Moré and Sorensen [213]) compute a direction of negative curvature, whenever it exists, to prevent the iteration from converging to nonminimizing stationary points. A direction of negative curvature  $p_-$  is one that satisfies  $p_-^T \nabla^2 f(x_k) p_- < 0$ . These algorithms generate a search direction by combining  $p_-$  with the steepest descent direction  $-\nabla f_k$ , often performing a curvilinear backtracking line search.

It is difficult to determine the relative contributions of the steepest descent and negative curvature directions. Because of this fact, the approach fell out of favor after the introduction of trust-region methods.

For a more thorough treatment of the modified Cholesky factorization see Gill, Murray, and Wright [130] or Dennis and Schnabel [92]. A modified Cholesky factorization based on Gershgorin disk estimates is described in Schnabel and Eskow [276]. The modified indefinite factorization is from Cheng and Higham [58].

Another strategy for implementing a line search Newton method when the Hessian contains negative eigenvalues is to compute a direction of negative curvature and use it to define the search direction (see Moré and Sorensen [213] and Goldfarb [132]).

Derivative-free line search algorithms include golden section and Fibonacci search. They share some of the features with the line search method given in this chapter. They typically store three trial points that determine an interval containing a one-dimensional minimizer. Golden section and Fibonacci differ in the way in which the trial step lengths are generated; see, for example, [79, 39].

Our discussion of interpolation follows Dennis and Schnabel [92], and the algorithm for finding a step length satisfying the strong Wolfe conditions can be found in Fletcher [101].

---

## EXERCISES

- 3.1 Program the steepest descent and Newton algorithms using the backtracking line search, Algorithm 3.1. Use them to minimize the Rosenbrock function (2.22). Set the initial step length  $\alpha_0 = 1$  and print the step length used by each method at each iteration. First try the initial point  $x_0 = (1.2, 1.2)^T$  and then the more difficult starting point  $x_0 = (-1.2, 1)^T$ .
- 3.2 Show that if  $0 < c_2 < c_1 < 1$ , there may be no step lengths that satisfy the Wolfe conditions.
- 3.3 Show that the one-dimensional minimizer of a strongly convex quadratic function is given by (3.55).
- 3.4 Show that the one-dimensional minimizer of a strongly convex quadratic function always satisfies the Goldstein conditions (3.11).
- 3.5 Prove that  $\|Bx\| \geq \|x\|/\|B^{-1}\|$  for any nonsingular matrix  $B$ . Use this fact to establish (3.19).
- 3.6 Consider the steepest descent method with exact line searches applied to the convex quadratic function (3.24). Using the properties given in this chapter, show that if the initial point is such that  $x_0 - x^*$  is parallel to an eigenvector of  $Q$ , then the steepest descent method will find the solution in one step.