# Continuous Optimization

· Unconstrained problem

$$\min \quad f(x) \qquad \qquad \underline{x} \in R^{N}$$

## Important cases

$$f \in C^{2}$$

$$f \in C^{1}$$

$$f \in C^{0}, \quad \text{convex}$$

$$f \in C^{0}$$

---

$\underline{f \in C^{2}}$:  Continuous, cont derivative, 2nd derivative

a  <u>Descent Direction</u> $\underline{d}$ has

the  property

$$f(\underline{x}_0 + \epsilon \underline{d}) \quad < \quad f(x_0) \qquad \text{for}$$

$$\epsilon < \Gamma_d$$

There are <u>numerous</u> descent directions

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

<u>this is gradient descent</u>

<u>issues :</u>
· how to choose $\varepsilon$

· <u>perhaps</u> interval halving
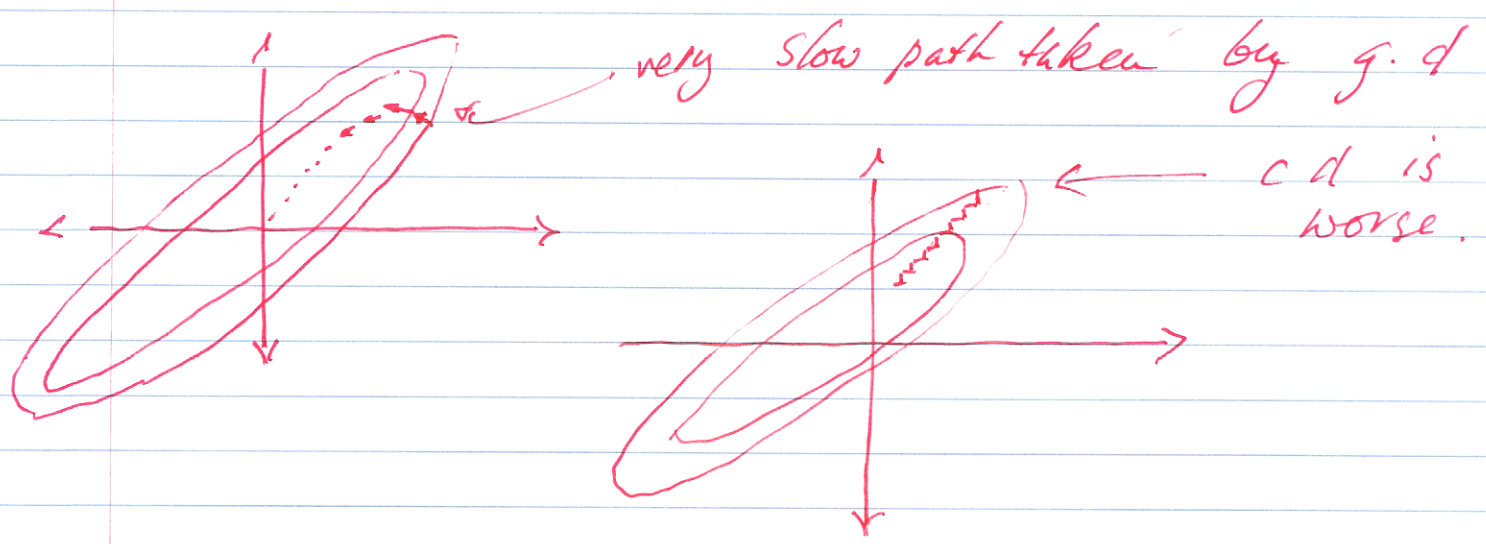
· more sophisticated machinery later

* 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 <u>coordinate descent</u>

* Both gradient, coordinate descent are naughty.

very slow path taken by g.d

c d is worse.

* $$H_f d_N = -\nabla f$$

<u>( Newtons method</u>

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

i.e. min $\nabla f \delta + \frac{1}{2} \delta' H_f \delta$

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

_Q_: Does Newton's method always give a Descent Direction?

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

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

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

— So we're ok if $H_f$ is positive <u>Definite</u>

_Q_: is $p$ a Descent Direction?

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

<u>Notice</u> we can obtain Descent Dirs from

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

$p$

iteration.

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

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

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

but

$B_k$ must be P.d for $p$ to be a descent dirn.

Example:  coordinate descent and EM

we have two parametric models

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

and ~~a noise model~~

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

and we observe $x_i$ from a mixture

$$p(x|\theta) = \mu_1 P_1 + \mu_2 P_2$$

write  CDLLH

$$\log\left[P(x_i, \delta_i \mid \theta)\right] = \sum_i \log\left[P(x_i \mid \delta_i, \theta)\, P(\delta_i)\right]$$

$$= \sum_i \left[-\delta_i\, g(x_i; \theta_1) + \delta_i \log \mu_i\right]$$

$$\overset{(1-\delta_i)}{\sum_i} \left[(1-\delta_i)\, g(x_i; \theta_2) + \log(1-\mu_i)\right]$$

$$= \mathcal{L}$$

But $\delta_i$ are unknown.

EM:  • start with $\theta^{(0)}\, \mu^{(0)}$

⌐→  • form $Q(\theta; \theta^{(i)}) = E_{\delta_i \mid \theta, \mu^{(i)}}\left[\mathcal{L}(\theta, \delta)\right]$

⌐·  • now form $\theta^{(i+1)} = \arg\max_\theta Q(\theta; \theta^{(i)})$

## EM for our example.

### E - step.

$$\bar{\delta}_i = 1 \cdot P(\delta_i = 1 \mid x_i, \theta)$$

$$= \frac{P(x_i \mid \delta_i = 1, \theta)\, P(\delta_i = 1 \mid \theta)}{P(x_i \mid \delta_i = 0, \theta)\, P(\delta_i = 0, \theta) + P(x_i \mid \delta_i = 1, \theta)\, P(\delta_i = 1 \mid \theta)}$$

$$= \frac{e^{-g_1(x_i, \theta^{(n)})} \cdot \mu^{(n)}}{e^{-g_2(x_i, \theta^{(n)})}(1 - \mu^{(n)}) + e^{-g_1(x_i, \theta^{(n)})} \cdot \mu^{(n)}}$$

### M - step:   Substitute + max

Now, consider

$$\mathcal{F}(\theta, \delta)$$

$$= C(\theta, \delta) + H(\delta)$$

↑ COLLH    ↑ entropy

$$H(\delta) = -\sum_i \left[ \delta_i \log \delta_i + (1 - \delta_i) \log(1 - \delta_i) \right]$$

Consider:
$$\nabla_\theta \mathcal{F}(\theta, \delta^{(n)}) = 0$$

↰ this is our old M step.

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

$$\frac{\partial \mathcal{F}}{\partial \delta_i} = -\left[ \log \delta_i - \log(1 - \delta_i) \right] + \left[ -\frac{\delta_i}{\delta_i} g_1 + \frac{\delta_i}{\delta_i} \log \mu^{(n)} \right]$$

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

i.e.

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

hence: <u>EM is coordinate ascent</u>

Q: why not do newtons method?

A: not sure, frankly.

H is <u>big</u>, but <u>sparse</u>

Issues:

- What to do with a descent dirn?
- 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 $\alpha$.

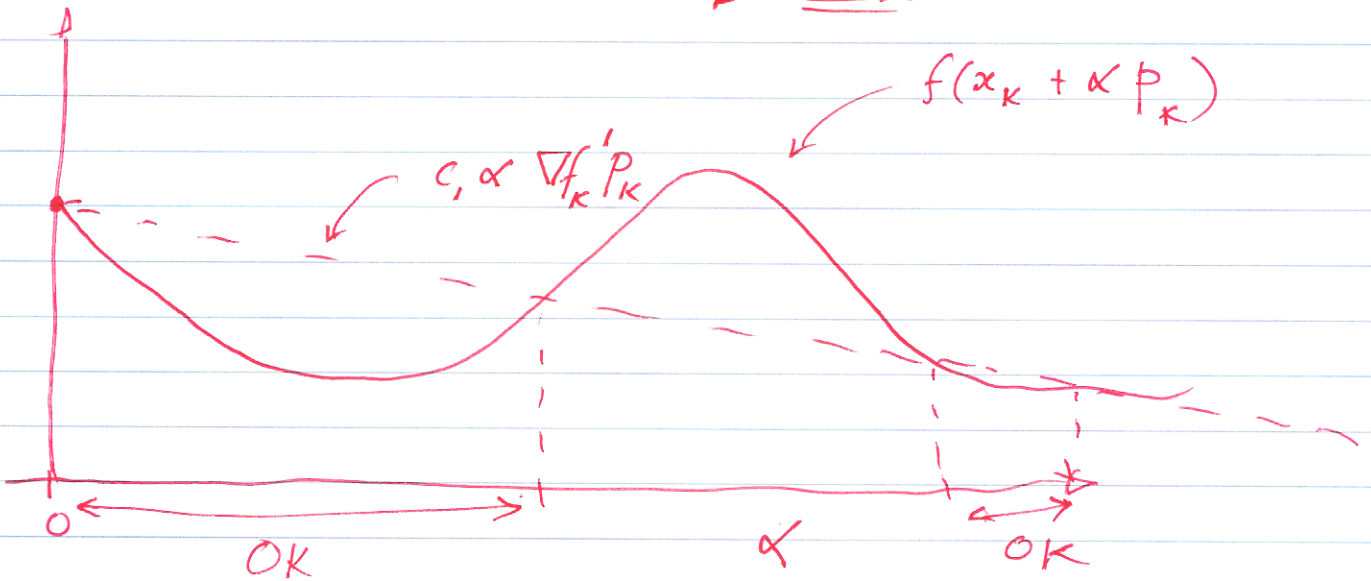consider $f(x_k + \alpha p_k)$     $\alpha > 0$

<u>what $\alpha$ are acceptable ?</u>

- 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 (typically $1ee^-4$) $\left[ \begin{array}{l} \underline{Armijo} \\ \underline{Wolfe} \end{array} \right\} Condition \right]$



$f(x_k + \alpha p_k)$

$C_1 \alpha \nabla f_k^T p_k$

$O$

$OK$          $\alpha$          $OK$

Sufficient Decrease is <u>not</u> enough

→ very small $\alpha$ are OK.

$$\cancel{\forall \alpha} \qquad \nabla f(x_k + \alpha_* p_k)^T p_k \geqslant c_2 \nabla f_k^T p_k$$

$$c_1 < c_2 < 1$$

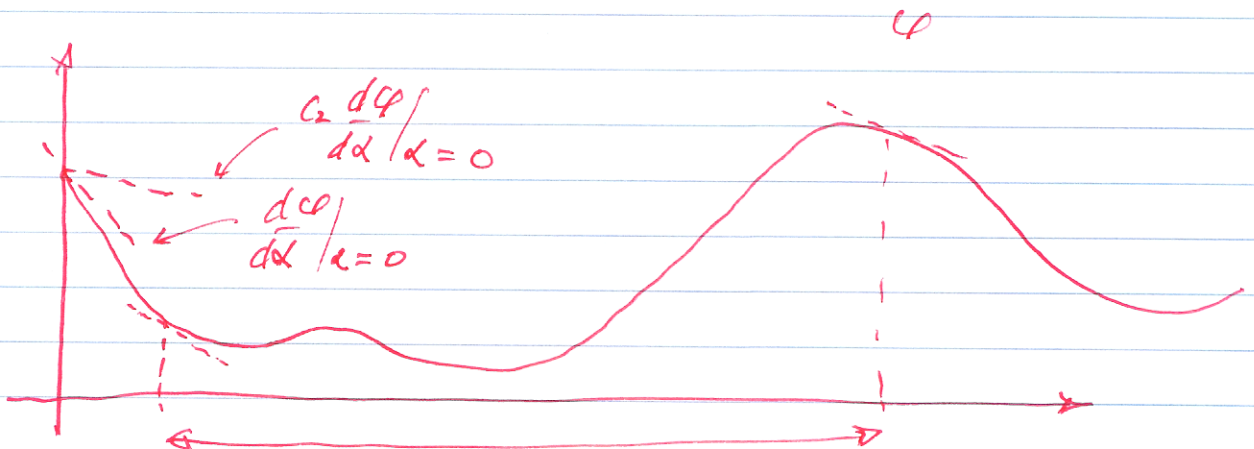<u>notice</u> :

write $\qquad \varphi(\alpha) = f(x_k + \alpha p_k)$

then $\qquad \dfrac{d\varphi}{d\alpha} = \nabla f(x_k + \alpha p_k)^1 p_k$

so condition is :

$$\dfrac{d\varphi}{d\alpha} \quad \cancel{\leqslant} \geqslant \quad c_2 \nabla f_k^T p_k = c_2 \dfrac{d\varphi}{d\kappa}\Big|_{\kappa=0}$$



$c_2 \dfrac{d\varphi}{d\alpha}\Big|_{\alpha=0}$

$\dfrac{d\varphi}{d\alpha}\Big|_{\alpha=0}$

$\varphi$

OK

Notice Sign of Slope!

$c_2$ is usually $0.4$ $\qquad$ (Newton)

$\qquad\qquad\qquad\qquad 0.1$ $\qquad$ (conj. grad.)

**Wolfe conditions**

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

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

**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 ]

$\qquad \alpha = \rho \alpha$

end

OK for newton; not as good for others.

## Newtons method with Hessian modification

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

so

$$H p_k = -\nabla f \quad \text{may not give}$$

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

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

$$\Rightarrow$$

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

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

1: <u>Add a multiple of Identity:</u>

choose $\beta > 0$
if $\min_i h_{ii} > 0$
$\qquad \tau_0 = 0$
else
$\qquad \tau_0 = - \min(h_{ii}) + \beta;$
end

for $k = \cdots$

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

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

$\left(\text{we are searching for } \tau I \text{ to}\right.$
make $H$ p.d.

## Cholesky:

for $j = 1 \cdots n$

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

$$d_j = c_{jj}$$

for $i = j+1 \cdots n$

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

$$\ell_{ij} = c_{ij} / d_j$$

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

$$d_j \geq \delta$$

$$\boxed{m_{ij} = \ell_{ij} \sqrt{d_j}} \quad \leq \quad \beta$$

Desirable for error control.

- <u>Improvements</u>

    - 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 $\alpha_0$

   $\rightarrow$ OK ; stop

   $\rightarrow$ Not OK ; there is an OK step
   in interval.

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

  - build quadratic interpolate

   $$\frac{\varphi(\alpha)}{q} = \left( \frac{\varphi(\alpha_0) - \varphi(0) - \alpha_0 \varphi'(0)}{\alpha_0^2} \right) \alpha^2$$

   $$+ \quad \varphi'(0) \alpha$$

   $$+ \quad \varphi(0)$$

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

$\longrightarrow$    $\alpha_1$ OK ;    stop

$\longrightarrow$    else    construct    <u>cubic</u>

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

$$\varphi(\alpha_1)$$

minimize ;    $\alpha_2$

$\longrightarrow \alpha_2$ OK     stop

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

two most recent $\alpha$

It can be shown that if $x_k \to x^*$ 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 $i > 1]$
            $\alpha_* \leftarrow$ **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 $\alpha_* \leftarrow$ **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
we can simply set $\alpha_{i+}$
important that the suc
a finite number of iter

We now specify
its input arguments is

  **(a)** the interval boun
       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

**Algorithm 3.6** (zoom
  **repeat**
        Interpolate (u
            a trial st
        Evaluate $\phi(\alpha_j$
        **if** $\phi(\alpha_j) > \phi($
            $\alpha_{hi} \leftarrow$
        **else**
            Evaluate
            **if** $|\phi'(\alpha_j$
                S
            **if** $\phi'(\alpha_j$
                $\alpha$
            $\alpha_{lo} \leftarrow \alpha$
  **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 $x_{\mathrm{lo}}$, then we set $\alpha_{\mathrm{lo}} \leftarrow \alpha_j$ to maintain condition (b). If this setting results in a violation of condition (c), we remedy the situation by setting $\alpha_{\mathrm{hi}}$ to the old value of $\alpha_{\mathrm{lo}}$. 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.