



UNITED NATIONS EDUCATIONAL, SCIENTIFIC AND CULTURAL ORGANIZATION
INTERNATIONAL ATOMIC ENERGY AGENCY
INTERNATIONAL CENTRE FOR THEORETICAL PHYSICS
I.C.T.P., P.O. BOX 586, 34100 TRIESTE, ITALY, CABLE: CENTRATOM TRIESTE



H4.SMR/994-5

**SPRING COLLEGES IN
COMPUTATIONAL PHYSICS**

19 May - 27 June 1997

**LINEAR ALGEBRA, MINIMIZATION, AND
EIGENVALUE EQUATIONS**

C. REBBI
Boston University
Department of Physics
590 Commonwealth Ave.
Boston, Massachusetts 01125
U.S.A.

*Preliminary notes from a course on Computational Physics -
copyright by Claudio Rebbi, 1995 - for distribution to participants only.*

Related problems:

- i) finding the zeros of several equations in several unknowns;
- ii) minimizing a function of several variables;
- iii) solving sets of linear equations.

$$\text{i) } f_i(x_j) = 0 \quad 1 \leq i, j \leq N.$$

Near the zeros we can expand

$$f_i(x_j) = f_i(x_j^{(0)}) + \left. \frac{\partial f_i}{\partial x_j} \right|_{x_j=x_j^{(0)}} (x_j - x_j^{(0)}) + \dots$$

\nearrow
implicit sum over j

$$\stackrel{\text{definition}}{=} A_{ij} x_j - b_i + \dots \equiv Ax - b \dots$$

compact notation:
 x, b are vectors
 A is a matrix

As in the algorithm of Newton + Raphson we take $x^{(0)}$ to be the first iterate and find a new zero by solving

$$Ax - b = 0,$$

i.e.

$$x^{(1)} = A^{-1} b.$$

Thus i) leads to iii).

ii) Minimization of $F(x_i)$.

At the minimum $f_i(x_j) = \frac{\partial F}{\partial x_i}(x_j) = 0$.
 ↓
 definition

This is a set of N equations in N unknowns, so that ii) leads to i).

Vice versa, i) can be reduced to ii) by defining

$$\tilde{F}(x_j) = \sum_i \varphi_i f_i^2(x_j) \quad (\alpha, \text{ with suitable weights } \varphi_i > 0)$$

$\tilde{F}(x_j) = \sum_i \varphi_i f_i^2(x_j)$ and demanding that

$\tilde{F}(x_j)$ is minimized (in this case $\min \tilde{F}(x_j) = 0$).

ii) \rightarrow iii) In the neighborhood of the minimum

the first derivatives $f_i(x_j) = \frac{\partial F}{\partial x_i}(x_j)$ can be

expanded as

$$f_i(x_j) = f_i(x_j^{(0)}) + \left. \frac{\partial^2 F}{\partial x_i \partial x_j} \right|_{x=x^{(0)}} (x_j - x_j^{(0)}) + \dots$$

↗
Hessian matrix

$$= A_{ij} x_j - b_i$$

and we can take as next iterate $x^{(1)} = A^{-1} b$.

Equivalently, we can think of expanding $F(x_i)$

directly :

$$\begin{aligned}
 F(x_i) = & F(x_i^{(0)}) + f_i(x_j^{(0)})(x_i - x_i^{(0)}) + \\
 & + \frac{1}{2} \left. \frac{\partial^2 F}{\partial x_i \partial x_j} \right|_{x=x^{(0)}} (x_i - x_i^{(0)})(x_j - x_j^{(0)}) + \\
 & + \dots
 \end{aligned}$$

In the neighborhood of the minimum the terms of higher order will be small and $F(x_i)$ will be well approximated by the quadratic form given by the first three terms in the expansion above.

Using the same notation as in the previous page, we can write this quadratic form as,

$$Q(x) = \frac{1}{2} A_{ij} x_i x_j - b_i x_i + c.$$

The minimum of Q is at x solving the equations

$$A x - b = 0$$

$$\text{i.e. } x = A^{-1} b.$$

Then we take as new iterate, etc..

Note: for a well posed problem of minimization the Hessian matrix is positive definite, i.e.

$$x_i A_{ij} x_j > 0 \text{ for all } x_i \neq 0.$$

In 1 dimension at a local minimum



$$f''(x^*) > 0$$

and, by continuity,

$f''(x) > 0$ also in a neighbourhood of

\underline{x}

The equivalent of the recent method case also be devised,

e.g., in 2 dimensions, imagine we wish to solve

and $f_1(x_1, x_2) = 0$

$$f_2(x_1, x_2) = 0.$$

Let us calculate f_1 in 3 points P_1, P_2, P_3 . We can then approximate $f_1(x_1, x_2)$ by a linear function

$$f_1(x_1, x_2) \approx Q_{11}x_1 + Q_{12}x_2 - b,$$

which takes the same values as f_1 over P_1, P_2, P_3 (through 3 points).

Similarly we calculate also f_2 at P_1, P_2, P_3 and approximate

$$f_2(x_1, x_2) \approx Q_{21}x_1 + Q_{22}x_2 - b.$$

As next iterate we use then the common zero of the

two planes, i.e. the solution of

$$Ax - b = 0.$$

Let us denote this solution by P_4 . We repeat
here the procedure on the basis of $P_1 P_2 P_3$
etc...

ii) \rightarrow iii) For a linear system of eqs.

$$Ax - b = 0$$

the solution minimizes the quadratic form

$$Q(x) = (x^T A^T - b^T) (Ax - b) =$$

$\uparrow \quad \uparrow$
let us drop the T from x, b

$$= x^T A^T A x - b^T A x - x^T A^T b - b^T b$$

The fact that $Q(x)$ is minimized
by the solution of $Ax - b = 0$ is obvious, and,
in any case, we can also demand

$$\frac{\partial}{\partial x_i} Q(x) = 0,$$

which gives $2A^T A x - 2A^T b = 0$.

For non-singular A this is equivalent to

$$Ax - b = 0.$$

Notice that the minimization of $Q(x)$ is

equivalent to the minimization of

$$\hat{Q}(x) = \frac{1}{2} x^T A^T A x - x^T A^T b.$$

If A is already symmetric and positive definite we can consider directly the quadratic form

$$Q = \frac{1}{2} x^T A x - x^T b$$

and minimize. Finding the minimum is equivalent to solving

$$Ax - b = 0.$$

\Leftarrow

Methods for solving a system of linear equations:

$$Ax = b.$$

The method of Gaussian elimination:

let A be given by

$$A = \begin{pmatrix} a_{11} & a_{12} & \dots \\ a_{21} & a_{22} & \dots \\ \dots & \dots & \dots \end{pmatrix}$$

and let us denote by r_1, r_2, \dots its rows,

i.e.

$$r_1 \equiv a_{11} \ a_{12} \ a_{13} \ \dots$$

$$r_2 \equiv a_{21} \ a_{22} \ a_{23} \ \dots$$

Replacing Γ_2 with $\Gamma_2 - \frac{a_{21}\Gamma_1}{a_{11}}$ the matrix becomes

$$A' = \begin{pmatrix} a_{11} & a_{12} & \dots \\ 0 & a_{22}' & \dots \\ & & \text{etc...} \end{pmatrix}$$

But this corresponds to

$$A \rightarrow A' = S_{12} A$$

with

$$S_{12} = \begin{pmatrix} 1 & 0 & 0 & \dots \\ -\frac{a_{21}}{a_{11}} & 1 & 0 & \dots \\ 0 & 0 & 1 & \dots \\ & & & \dots \end{pmatrix}$$

Obviously $Ax = b$ is equivalent

$$\text{to } S_{12} A x = S_{12} b.$$

With many similar replacements (Gaussian elimination) the system becomes

$$\hat{A}x = \hat{b}$$

with \hat{A} upper triangular.

(Notice, \hat{A} and A also have the same determinant,
 (but $\det \hat{A} = \hat{a}_{11} \hat{a}_{22} \dots \hat{a}_{nn}$, so that the method
 of Gaussian elimination also provides a collapse -
 (ratio of the determinant)).

With \hat{A} upper triangular

$$\hat{A}^T x = \hat{b} \Rightarrow \text{ solved immediately}$$

by backward substitution

\uparrow

$$② \text{ gives } x_{n-1} \quad \hat{a}_{n-1, n-1} x_{n-1} + \hat{a}_{n-1, n} x_n = b_{n-1},$$

$$① \text{ gives } x_n \quad \hat{a}_{nn} x_n = b_n$$

Pivoting is crucial: one should not use a_{ii} directly, but rather should transpose rows and columns so as to bring into the a_{ii} position the matrix element with the largest absolute value, etc..

Operations' count: $O(N)$ operations for every \hat{a} introduced in $\hat{A}' \hat{A}''$ etc., hence $O(N^3)$ operations to bring A to upper triangular form;
 $O(N^2)$ operations for the backward substitution.

Solving $Ax = b$, finding A^{-1} and $\det A$,
 (as well as the eigenvalues and eigenvectors of A)
 are typical operations for which there exist good
 library subroutines.

A typical library subroutine might be of
 the form

DLAT (N, V (A, B, M, N, M1, N1, DET, IERR))

and might produce the solution to the equations

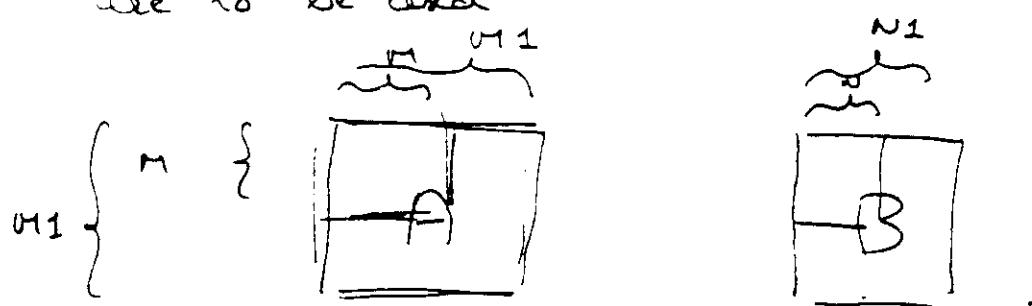
$$A X = B,$$

where A is an $M \times N$ matrix and X, B
 are $N \times 1$ vectors (this is equivalent to
 solving $Ax = b$ for N different vectors b).

If you return, the diagonal elements of B
 might be overwritten with the solution X
 (note, if an entry $B \approx 0$, then on return it
 will contain (A^{-1})), DET will contain the
 determinant of A and IERR some error code,
 i.e. something went wrong.

M1 and N1 might specify some physical
 dimensions, into which A and B are
 embedded. I.e., it might be that,
 for convenience, A is dimensioned as a

$M_1 \times M_1$ matrix, B as a $M_1 \times N_1$ matrix,
of which only the upper $M \times M$, $M \times N$ elements
are to be used



Because of the way the matrices are stored,

both M and M_1 ,
must be reversed.

Want N_1

Special routines will be used for more efficient
solutions in special cases, e.g. A could be
symmetric etc -

L U decomposition:

notice, the gaussian elimination procedure
implicitly does implement an LU decomposition.

A is transformed by

$$A \rightarrow S_{12} A \rightarrow S_{13} S_{12} A \rightarrow \dots \text{ until}$$

$$S_{N-1N} \dots S_{12} A = U \quad \text{upper triangular matrix}$$

Thus

$$A = S_{12}^{-1} \dots S_{n-1n}^{-1} U .$$

But

$$S_{12}^{-1} = \begin{pmatrix} 1 & 0 & 0 & \dots \\ \frac{a_{21}}{a_{11}} & 1 & 0 & \dots \\ 0 & 0 & 1 & \dots \\ \vdots & \vdots & \vdots & \ddots \end{pmatrix}$$

etc.

so that $A = L U$ with L lower triangular;
Moreover all diag. elements of L are 1.

The decomposition can be done directly:

we proceed

$$\begin{pmatrix} 1 & 0 & 0 & \dots & | & u_{11} & u_{12} & u_{13} & \dots \\ l_{21} & 1 & 0 & \dots & | & 0 & u_{22} & u_{23} & \dots \\ l_{31} & l_{32} & 1 & \dots & | & 0 & 0 & u_{33} & \dots \\ \vdots & \vdots & \vdots & \ddots & | & \vdots & \vdots & \vdots & \ddots \end{pmatrix} = \begin{pmatrix} a_{11} & a_{12} & a_{13} & \dots \\ a_{21} & a_{22} & a_{23} & \dots \\ a_{31} & a_{32} & a_{33} & \dots \\ \vdots & \vdots & \vdots & \ddots \end{pmatrix}$$

The multiplication gives

$$\begin{array}{cccccc} | & u_{11} & u_{12} & u_{13} & \dots \\ | & l_{21} u_{11} & l_{21} u_{12} + u_{21} & l_{21} u_{13} + u_{23} & \dots \\ | & l_{31} u_{11} & l_{31} u_{12} + l_{32} u_{21} & l_{31} u_{13} + l_{32} u_{23} + u_{31} & \text{etc.} \dots \\ | & & & & & \end{array}$$

Equating the first row to $a_{11} a_{12} a_{13} \dots$
 immediately fixes $u_{11} u_{12} u_{13} \dots$

The second row fixes l_{21} first, then $u_{22} u_{23} \dots$

The third row fixes l_{31} and l_{32} first, then $u_{33} u_{34} \dots$
 etc.

With $A = LU$ the linear system of equations
 $Ax = b$ can be recast into

$$\begin{aligned} Ux &= y \\ Ly &= b \end{aligned} \quad (Ly = LUx = b)$$

$Ly = b$ is now solved immediately by forward substitution
 and $Ux = y$ is solved by backward substitution.

One should rearrange rows and columns so that
 u_{11} is the largest element, in abs. value, u_{22} the
 largest in the remaining $(N-1) \times (N-1)$ square etc..

This does not require actual rearrangements,
 but simply a suitable indexing. The arrays
 $iR(N)$ and $iC(N)$ will contain the
 indices of the rows and columns in the order
 in which they would be rearranged if we had
 actually performed the substitutions.

For example, with

$$\begin{pmatrix} 4 & 9 & 2 \\ 2 & 8 & 10 \\ 1 & 7 & 3 \end{pmatrix}$$

ir and ic would contain $ir = (2 \ 1 \ 3)$ $ic = (3 \ 2 \ 1)$

because 10 ($r=3 c=3$) is the largest element.

9 ($r=1 c=2$) is the largest in the remaining

square

$$\begin{array}{cc|c} 4 & 9 & | \\ \hline 1 & 7 & | \end{array}$$

etc...

We then use an indirect addressing of rows
and columns

DO $ia = 1, N$

DO $ja = 1, N$

$ib = ir(ia)$

$jb = ic(ja)$

use ib jb as indices

An annealing algorithm for symmetric matrices.

Given a symmetric M find M^{-1} .

Write M as

$$M = \begin{pmatrix} A & b \\ b^T & c \end{pmatrix}$$

$M = N \times N$ matrix $A = (N-1) \times (N-1)$ $b = (N-1) \times 1$ $c = \text{constant}$

$$M = \left(\begin{array}{c|c} \overline{A} & \overline{b} \\ \hline b^T & c \end{array} \right)$$

Write M^{-1} in the form $\begin{pmatrix} X & y \\ y^T & z \end{pmatrix}$,

and solve

$$\begin{pmatrix} X & y \\ y^T & z \end{pmatrix} \begin{pmatrix} A & b \\ b^T & c \end{pmatrix} = \begin{pmatrix} I & 0 \\ 0 & 1 \end{pmatrix}.$$

This gives

- 1) $X A + y b^T = \underline{\underline{I}} ,$
- 2) $X b + c y = 0 ,$
- 3) $y^T A + z b^T = 0 , \text{ or } A y + z b = 0$
- 4) $y^T b + z c = 1 \quad (\text{or } \bar{y} \cdot \bar{b} + z c = 0).$

3) gives $y = -z A^{-1} b .$

With this 1) becomes $X A - z A^{-1} b b^T = \underline{\underline{I}} ,$

which gives

$$\begin{aligned} X &= \underline{\underline{A}}^{-1} - z (A^{-1} b) (b^T A^{-1}) = \\ &= \underline{\underline{A}}^{-1} + \frac{y y^T}{z} . \end{aligned}$$

4) fixes z , since it gives $-z b^T A^{-1} b + z c = 1$

$$\text{i.e. } z = \frac{1}{c - b^T A^{-1} b} .$$

2) is also solved : $X b + c y = (\underline{\underline{A}}^{-1} + \frac{y y^T}{z}) b + c y =$
 $= A^{-1} b - y (b^T A^{-1} b) + c y =$
 $= -\frac{y}{z} - y (b^T A^{-1} b) + c y =$
 $= - (c - b^T A^{-1} b) y + (c - b^T A^{-1} b) y = 0 .$

Thus, given $M = \begin{pmatrix} A & b \\ b^T & c \end{pmatrix}$ if we knew A^{-1}

we can calculate

$$M^{-1} = \begin{pmatrix} A^{-1} + \frac{yy^T}{z} & y \\ y^T & z \end{pmatrix},$$

with $v = A^{-1}b$

$$z = \frac{1}{c - b^T v}$$

$$y = -zv.$$

Notice that we do not need to know A .

Thus we can invert M by induction, working out from the inverse of its first element m_{11} .

$$\begin{pmatrix} m_{11} & m_{12} & \dots \\ m_{21} & m_{22} & \dots \\ \vdots & \vdots & \ddots \end{pmatrix} \Rightarrow \begin{pmatrix} \underline{m_{11}} & m_{12} & \dots \\ m_{21} & m_{22} & \dots \\ \vdots & \vdots & \ddots \end{pmatrix} \Rightarrow$$

$$= \begin{pmatrix} (m_{11} \ m_{12})^{-1} & m_{13} \\ m_{21} \ m_{22} & m_{23} \\ m_{31} \ m_{32} \ m_{33} \end{pmatrix} \Rightarrow \text{etc.} \dots$$

For example, write a 2×2 matrix

$$\text{Let } \begin{pmatrix} a & b \\ b & c \end{pmatrix} \text{ we get } \begin{pmatrix} a^{-1} & b \\ b & c \end{pmatrix} \text{ and}$$

$$u = a^{-1}b \quad z = \frac{1}{c - a^{-1}b^2} \quad y = -\frac{a^{-1}b}{c - a^{-1}b^2}$$

$$z = \frac{a}{ac - b^2} \quad y = -\frac{b}{ac - b^2}$$

$$M^{-1} = \begin{pmatrix} a^{-1} + \frac{b^2(c-a^{-1}b^2)}{(ac-b^2)a} & -\frac{b}{ac-b^2} \\ -\frac{b}{ac-b^2} & \frac{a}{ac-b^2} \end{pmatrix} =$$

$$= \begin{pmatrix} \frac{1}{a} & \frac{ac + b^2 + b^2}{ac - b^2} = \frac{c}{ac - b^2} & -\frac{b}{ac - b^2} \\ -\frac{b}{ac - b^2} & \frac{a}{ac - b^2} \end{pmatrix}$$

The analytic calculation looks rather complicated, but the iterative procedure is very simple.

With very sparse matrices, iterative techniques for the solution of a system of linear eqs.

$$Ax = b$$

may be more efficient than the methods illustrated above.

If A is symmetric and positive definite, solving the eq. $Ax=b$ is equivalent to finding the minimum of

$$Q(x) = \frac{1}{2} x^T A x - b^T x.$$

If A is not symmetric and positive definite we may replace the original system of eqs. with

$$A^T A x = A^T b,$$

or with $A^+ A x = A^+ b$ ($A^+ = \text{conjugate of } A^\dagger$)

if A is complex.

$$M = A^T A \quad \text{or} \quad M = A^+ A \Rightarrow \text{these}$$

symmetric and positive definite (a hermitian and p.d.).

Let us study the problem of minimizing

$$Q(x) = \frac{1}{2} x^T A x - b^T x.$$

Search along the steepest descent.

Let us start from a point x_0 . We can calculate

$$g_0 = -\text{grad } Q \Big|_{x=x_0} = -Ax_0 + b.$$

We may now replace x_0 with

$$x_1 = x_0 + \beta g_0,$$

where β is a number which we will adjust so that the minimum of Q along the line starting at x_0 and directed along g_0 occurs at x_1 .

Later we will have to move along a direction $h \neq g_0$, so let us say that we move along the line

$$x_1 = x_0 + \beta h_0,$$

with $h_0 = g_0$ to begin with.

Let us find β :

at x_1 the negative gradient will be

$$g_1 = -Ax_1 + b = -Ax_0 - \beta Ah_0 + b = g_0 - \beta Ah_0.$$

But, if Q is minimised along the line, then g_1 must be orthogonal to the direction of the line, i.e. h_0 : $g_1 \cdot h_0 = 0$.

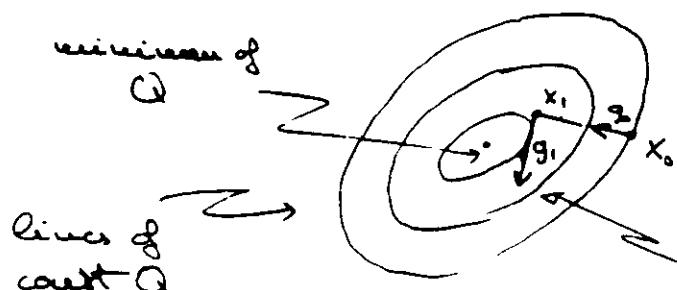
This fixes β :

$$g_0 h_0 = 0 \Rightarrow g_0 h_0 - \beta h_0 A h_0 = 0 \\ \text{i.e. } \beta = \frac{g_0 h_0}{h_0 A h_0}.$$

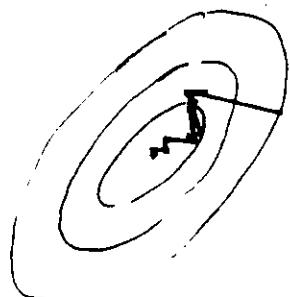
Do we now move in direction g_1 ?

No, g_1 is orthogonal to the original g_0 and the direction of steepest descent (i.e. g_0) is not the best.

See it in 2 dimensions:



g_1 is orthogonal to g_0 and clearly does not represent the best direction to approach the minimum.



If we kept on going along the steepest descent we would approach the solution along a path like the one on the left, with very poor convergence indeed.

Patent, we move along a direction h_1 , such that the new gradient is orthogonal to h_0 after the move.

Find the new direction of minimization:

$$\text{set } h_1 = g_1 + \alpha_0 h_0 \leftarrow (\text{correction term})$$

We will move along h_1 to

$$x_2 = x_1 + \beta_1 h_1$$

and, at the minimizer, the new gradient will be

$$g_2 = g_1 - \beta_1 A h_1$$

β_1 is fixed by $g_2 h_1 = 0$, i.e.

$$\beta_1 = \frac{h_1 g_1}{h_1 A h_1}$$

If we also demand

$$g_2 h_0 = 0$$

this gives

$$0 = g_1 h_0 - \beta_1 h_0 A h_1$$

\parallel
0 from previous step

$$\text{i.e. } h_0 A h_1 = 0$$

This fixes h_1 , indeed we get

$$h_0 A g_i + \alpha_0 h_0 A h_0 = 0$$

or $\alpha_0 = - \frac{g_i A h_0}{h_0 A h_0}$

These conditions will be maintained throughout all iterations, namely

from the given h_i , g_i we calculate the new minimum along h_i :

$$x_{i+1} = x_i + \beta_i h_i$$

by decreasing $g_{i+1} h_i = 0$,

i.e., since $g_{i+1} = g_i - \beta_i A h_i$

$$\beta_i = \frac{h_i g_i}{h_i A h_i}$$

The previous construction of h_i , determined by the condition

$$h_{i-1} A h_i = 0$$

implies that $g_{i+1} h_{i-1} = 0$ as well (since $g_i h_{i-1} = 0$).

We will thus define a new direction of minimization

$$h_{i+1} = g_{i+1} + \alpha_i h_i$$

such that $h_{i+1} A h_i = 0$, which gives

$$\alpha_i = - \frac{g_{i+1} A h_i}{h_i A h_i}$$

This algorithm is referred to as the method of conjugate gradients.

By construction, it gives

$$A) \quad g_i \cdot h_j = 0 \quad \text{for } j = i-1, i-2, \dots$$

$$B) \quad h_i \cdot A h_j = 0 \quad \text{for } j = i-1.$$

It also gives

$$C) \quad g_i \cdot g_j = 0 \quad \text{for } j = i-1,$$

since from $h_j = g_j + \alpha_{j-1} h_{j-1}$ it follows that
 g_j is a linear combination of h_j and h_{j-1} and
 $g_i \cdot h_{i-1} = g_i \cdot h_{i-2} = 0$.

But one can prove that eqs. A) B) C) are actually true for all $j < i$.

By induction: assume it proved up to i ,
prove it then for $i+1$.

$$A) \quad g_{i+1} \cdot h_j = 0 ?$$

This is true by construction for $j = i$,
for $j < i$ we use

$$g_{i+1} = g_i - \beta_i A h_i$$

and $g_i \cdot h_j = 0$, $h_i \cdot A h_j = 0$ which are both
true for $j < i$ by the inductive hypothesis.

$$C) \quad g_{i+1}g_i = 0 ?$$

From $h_j = g_j + \alpha_{j-1}h_{j-1}$ we get

$$g_{i+1}g_j = g_{i+1}(h_j - \alpha_{j-1}h_{j-1}) = 0 \text{ because A)}$$

has now been proved true up to $i+1$.

$$B) \quad h_{i+1}Ah_i = 0 ?$$

This is true by construction for $j=i$.

For $j < i$ we use $h_{i+1} = g_{i+1} + \alpha_i h_i$ and
the identity to prove becomes

$$(g_{i+1} + \alpha_i h_i) Ah_j = 0 \text{ for } j < i.$$

But $h_i Ah_j = 0$ for $j < i$ because of the inductive hypothesis; so we must prove

$$g_{i+1}Ah_j = 0 \text{ for } j < i.$$

From

$$g_{j+1} = g_j - \beta_j Ah_j \text{ we get}$$

$$g_{i+1}(Ah_j) = g_{i+1} \frac{g_j - g_{j+1}}{\beta_j}$$

and, with $j < i$, $j+1 \leq i$ the r.h.s. is 0 on account of C.

Note that the orthogonality and conjugate orthogonality relations

$$g_i \cdot g_j = 0 \quad i \neq j$$

$$h_i \cdot A h_j = 0 \quad i \neq j$$

imply that, if the dimensionality of the space is n , the algorithm must terminate in at most n steps (i.e., starting from g_0 and h_0 , g_n and h_n must turn out to vanish, so that x_n is the required minimum).

In practice, the algorithm converges to a high degree of accuracy in a number of steps much smaller than n . ^{frequently}

From A) B) and C) more equally conceivable relations follow, e.g.

$$g_i = h_i - \alpha_{i-1} h_{i-1} \text{ gives } h_j \cdot A g_i = 0 \text{ for } j > i;$$

$$\cdot A h_i = \frac{g_i - g_{i-1}}{\beta_i} \text{ gives } g_j \cdot A h_i = 0 \text{ for } j > i+1;$$

$$\beta_i = \frac{h_i \cdot g_i}{h_i \cdot A h_i} = \frac{(g_i + \alpha_{i-1} h_{i-1}) g_i}{(g_i + \alpha_{i-1} h_{i-1}) A h_i} - \frac{g_i^2}{g_i \cdot A h_i},$$

and, especially,

$$\begin{aligned}\alpha_i &= -\frac{g_{i+1}^T A h_i}{h_i^T A h_i} = \frac{g_{i+1}^T (g_{i+1} - g_i)}{\beta_i (h_i^T A h_i)} = \text{ (with } \beta_i := \frac{h_i^T g_i}{h_i^T A h_i}) \\ &= \frac{g_{i+1}^T (g_{i+1} - g_i)}{h_i^T g_i} = \frac{g_{i+1}^T (g_{i+1} - g_i)}{(g_i + \alpha_{i-1} h_{i-1}) g_i} = \text{ (because of A).} \\ &= \frac{g_{i+1}^T (g_{i+1} - g_i)}{g_i^2} = \text{ (because of c)} = \frac{g_{i+1}^2}{g_i^2}\end{aligned}$$

with no reference to the matrix A.

In this form the algorithm generalizes to the non-linear case.

Let $F(x)$ be any function

we start from x_0 and calculate $g_0 = -\text{grad } F|_{x_0}$

We select a direction h_0 , which the first time will coincide with g_0 .

$$h_0 = g_0.$$

We move along the line

$$x = x_0 + \beta h_0$$

and we look for the minimum of F along this line

*)

Let the next minimum be at x_i (now, of course, $i=1$, but by using a generic index i we can calculate the basic iterative cycle). We calculate the (negative) next gradient $g_i = -\text{grad } F|_{x_i}$

and the parameter

$$\alpha_{i-1} = \frac{g_i^2}{g_{i-1}^2}$$

With this we define the direction for the next search

$$h_i = g_i + \alpha_{i-1} h_{i-1}$$

We find the minimum along this line and go to step *).

To the extent that $F(x)$ is well approximated by a quadratic form $Q(x)$ this procedure will converge to the minimum of F . Convergence is not guaranteed, but, by construction, the values of the $F(x_i)$ will form a non-increasing (general decreasing) sequence. If the sequence does not converge to the required degree of approximation, one can always restart the algorithm from the final value found for x . This will typically be closer to the minimum and the approximation of F with a quadratic form will be better.

In ^{very} many cases the algorithm works in a quite satisfactory manner.

Calculation of eigenvalues and eigenvectors of a symmetric (or hermitian) matrix.

- 1) bring the matrix to tridiagonal form

$$\hat{M} = \begin{pmatrix} a, b, 0 & 0 \\ b, a, b, 0 \\ 0 & b, a, \ddots \\ 0 & 0 & \ddots \end{pmatrix},$$

- 2) calculate the determinant

$$\det(\hat{M} - \lambda I),$$

- 3) find the zeros (eigenvalues),

- 4) calculate the eigenvectors.

Reduction to tridiagonal form: Householder method

$$M = \begin{pmatrix} a & \vec{v}^\top \\ \vec{v} & A \end{pmatrix}$$

dim M = n
 $a = M_{11}$
 \vec{v} = n-1 component vector

$$A = (n-1) \times (n-1) \text{ matrix}.$$

Consider the matrix

$$U = \begin{pmatrix} 1 & \vec{\phi}^T \\ \vec{\phi} & P \end{pmatrix}$$

with $P = \mathbb{I} - 2\vec{u}\vec{u}^T$ and $\vec{u}^T\vec{u} = 1$.

U is symmetric and orthogonal. The symmetry is obvious. Then

$$\begin{aligned} P^T P &= P^2 = (\mathbb{I} - 2\vec{u}\vec{u}^T)(\mathbb{I} - 2\vec{u}\vec{u}^T) = \\ &= \mathbb{I} - 4\vec{u}\vec{u}^T + 4\vec{u}(\vec{u}^T\vec{u})\vec{u}^T = \mathbb{I}, \end{aligned}$$

so that U is indeed orthogonal.

The similarity transformation

$$M \rightarrow M' = U^T M U = U M U \text{ gives}$$

$$M' = \begin{pmatrix} a & (P\vec{v})^T \\ (P\vec{v}) & PAP \end{pmatrix}.$$

We fix now \vec{u} so that $P\vec{v} = \vec{\omega}$, where $\vec{\omega}$ is of the form $(\omega, \emptyset, \emptyset, \emptyset - \emptyset)$.

Then M' will be of the form

$$\left(\begin{array}{ccccc} m_{11} & m_{12} & 0 & 0 & \dots \\ m_{21} & \overline{m_{22}} & & & \\ 0 & & & & \\ 0 & & & & \\ \vdots & & & & \\ \end{array} \right) \quad (\text{with } m_{11} = \alpha \\ m_{12} = m_{21} = \omega)$$

$\underbrace{\qquad\qquad\qquad}_{M''}$

The method can then be applied to the $n-1 \times n-1$ matrix M'' etc.. and in $n-2$ steps the whole matrix will be brought to tridiagonal form.

We must still determine \vec{u} . We have

$$\vec{w} = P \vec{v} = \vec{v} - 2\vec{u}(\vec{u}^T \vec{v}). \quad *)$$

The orthogonality of P (or a direct calculation from the above equation*) gives $\vec{w}^T \vec{w} = \vec{v}^T \vec{v}$, hence

$$w = \pm |\vec{v}|. \quad \text{Moreover, we also have (from *)}$$

$$\vec{v}^T \vec{w} = \vec{v}^T \vec{v} - 2(\vec{v}^T \vec{u})(\vec{u}^T \vec{v}), \\ \pm v_1 |\vec{v}| = |\vec{v}|^2 - 2(\vec{u}^T \vec{v})^2.$$

This allows us to solve for $\vec{u}^T \vec{v}$. We choose the sign

that under the r.h.s. the largest need fixed $\vec{u}^T \cdot \vec{v}$
place

$$(\vec{u}^T \cdot \vec{v})^2 = \frac{|\vec{v}|^2 + |\vec{u}_1 \cdot \vec{v}|^2}{2}$$

We can now finally solve for \vec{u}

$$\vec{u} = \frac{\vec{v} - \vec{\omega}}{2(\vec{u}^T \cdot \vec{v})} = \frac{1}{\sqrt{2(|\vec{v}|^2 - |\vec{u}_1 \cdot \vec{v}|^2)}} (\vec{u}_1 \neq |\vec{v}| \vec{1}, \vec{u}_2, \vec{u}_3, \dots)$$

L

Assume now that M is tridiagonal.

We next calculate

$$\det(M - \lambda I).$$

$M - \lambda I$ is also tridiagonal.

Consider then a tridiagonal symmetric matrix
 $T \equiv T_n$. Let T_1, T_2, \dots, T_n be the matrices
 formed by the first 1, 2, ..., n rows and
 columns:

$$T_3 \left\{ T_2 \left\{ \begin{array}{|c|} \hline T_1 = t_{11} & t_{12} = 0 \\ \hline t_{12} & t_{22} = t_{23} \\ \hline 0 & t_{23} = t_{33} \\ \hline \end{array} \right\} \right\} = t_{33}.$$

We have $\text{Det } T_i = (\text{Det } T_{i-1}) t_{ii} - (\text{Det } T_{i-2}) t_{ii}^2$

With $\text{Det } T_1 = t_{11}$, $\text{Det } T_0 = 1$ this allows us to calculate $\text{Det } T_n$ by iteration. The procedure is very fast, so $\text{Det } (M - \lambda I)$ can be calculated quickly for every given value of λ , and this can be used in a program that looks for the zeros of the determinant, i.e. the eigenvalues.

Once the eigenvalues are found, the eigenvectors can be determined by solving $n-1$ among the n eqs

$$M \vec{v} = \lambda \vec{v},$$

considering $n-1$ components of \vec{v} as unknowns and one component as given.

The eigenvector \vec{v} corresponding to a definite eigenvalue λ can also be found by solving the equations

$$(M - \lambda \mathbb{I} + \epsilon \mathbb{I}) \vec{v} = \vec{v}_0,$$

where \vec{v}_0 is an arbitrary vector (not orthogonal to \vec{v}_λ) and ϵ is a small number. The solution of this equation enhances the component of \vec{v}_0 along the direction of the eigenvector \vec{v}_λ .

Indeed, the solution

$$\vec{v} = (M - \lambda \mathbb{I} + \epsilon \mathbb{I})^{-1} \vec{v}_0,$$

in the basis formed by the eigenvectors themselves takes the form

$$v_{\lambda'} = \frac{v_{0\lambda'}}{\lambda' - \lambda + \epsilon}$$

$$\text{so that, by way recall } \epsilon, \quad v_\lambda = \frac{v_{0\lambda}}{\epsilon}$$

dominates over all the other $v_{\lambda'}$.

The above procedure can be repeated a few times to enhance the component along \vec{v}_λ even more.

Reduction to tridiagonal form: the Lanczos algorithm.

Let us call the matrix to be brought to tridiagonal form H . It is a real symmetric or complex hermitian matrix.

We construct a basis of vectors $|4_1\rangle |4_2\rangle \dots$ such that H is tridiagonal in that basis.

(It is convenient to use Dirac's notation:

$$\vec{\psi}_1 = |4_1\rangle, \vec{\psi}_1^+ = \langle 4_1|, \vec{\psi}_1^+ H \vec{\psi}_2 = \langle 4_1 | H | 4_2 \rangle \text{ etc}$$

We choose $|4_1\rangle$ arbitrarily (or by some motivation given) normalized so that $\langle 4_1 | 4_1 \rangle = 1$.

$|4_1\rangle$ and $H|4_1\rangle$ in general span a 2 dimensional subspace (otherwise, if $H|4_1\rangle = c|4_1\rangle, |4_1\rangle$ is an eigenvector, we have solved part of our problem, and can proceed to another guess $|4_2\rangle$).

We define $|4_2\rangle$ as a linear combination of $H|4_1\rangle$ and $|4_1\rangle$ such that $|4_2\rangle$ is also normalized to 1 and that it is orthogonal to $|4_1\rangle$.

We write



$$|\psi_2\rangle = c_2 (H|\psi_1\rangle - \alpha_{11}|\psi_1\rangle).$$

Then $\langle\psi_1|\psi_2\rangle = 0$ gives $\alpha_{11} = \langle\psi_1|H|\psi_1\rangle$
and c_2 can be easily determined

$$(c_2 = [\langle\psi_1|H|H|\psi_1\rangle - \alpha_{11}^2]^{-\frac{1}{2}}).$$

We now determine $|\psi_3\rangle$ as a linear combination of $H|\psi_2\rangle$, $|\psi_2\rangle$ and $|\psi_1\rangle$ such that $\langle\psi_3|\psi_2\rangle = 0$,
 $\langle\psi_3|\psi_1\rangle = 0$ and $\langle\psi_3|\psi_3\rangle = 1$:

$$|\psi_3\rangle = c_3 (H|\psi_2\rangle - \alpha_{22}|\psi_2\rangle - \alpha_{12}|\psi_1\rangle).$$

Our conditions give $\alpha_{22} = \langle\psi_2|H|\psi_2\rangle$
 $\alpha_{12} = \langle\psi_1|H|\psi_2\rangle$ and c_3 again
immediately follows.

In general, we construct $|\psi_{n+1}\rangle$ from $|\psi_n\rangle$ and $|\psi_{n-1}\rangle$ as

$$|\psi_{n+1}\rangle = c_{n+1} (H|\psi_n\rangle - \alpha_{nn}|\psi_n\rangle - \alpha_{n-1,n}|\psi_{n-1}\rangle)$$

in such a way that

$$\langle\psi_{n+1}|\psi_n\rangle = \langle\psi_{n+1}|\psi_{n-1}\rangle = 0, \quad \langle\psi_{n+1}|\psi_{n+1}\rangle = 1.$$

By construction, every $|4_n\rangle$ is orthogonal to $|4_{n-1}\rangle$ and $|4_{n-2}\rangle$. We can prove however that $|4_n\rangle$ is orthogonal to all $|4_i\rangle$ with $i < n$.

We proceed by induction. We assume that the property holds true up to $|4_n\rangle$ and prove it true for $|4_{n+1}\rangle$. We only need to prove that

$$\langle 4_i | 4_{n+1} \rangle = 0 \quad \text{for } i < n+1,$$

since for $i = n+1$ and $i = n$ the equality is satisfied by construction.

With $i < n+1$ we have $\langle 4_i | 4_n \rangle = \langle 4_i | 4_{n-1} \rangle = 0$ by the inductive hypothesis, so that

$$\begin{aligned} \langle 4_i | 4_{n+1} \rangle &= \langle 4_i | (c_{n+1}(H|4_n\rangle - a_{n+1}|4_n\rangle - b_{n+1}|4_{n-1}\rangle) \\ &= c_{n+1} \langle 4_i | H | 4_n \rangle. \end{aligned}$$

But from $|4_{i+1}\rangle = c_{n+1}(H|4_i\rangle - a_{n+1}|4_i\rangle - b_{n+1}|4_{i-1}\rangle)$ it follows that $H|4_i\rangle$ is a linear combination of $|4_{i+1}\rangle, |4_i\rangle, |4_{i-1}\rangle$. Since $i < n+1$, all these vectors are orthogonal to $|4_n\rangle$ by the inductive hypothesis. It follows that

$$\langle 4_i | H | 4_n \rangle = 0 \quad \text{and} \quad \langle 4_i | 4_{n+1} \rangle = 0.$$

Q.E.D.

Thus, by the Lanczos algorithm, we can construct a basis $|4_1\rangle, |4_2\rangle \dots$ of orthonormal vectors. But from the relation

$$|4_{n+1}\rangle = C_{n+1}(H|4_n\rangle - Q_{n+1}|4_n\rangle - Q_{n+1}^*|4_{n+1}\rangle),$$

i.e.

$$H|4_n\rangle = Q_{nn}|4_n\rangle + Q_{n+1,n}|4_{n+1}\rangle - \frac{|4_{n+1}\rangle}{C_{n+1}},$$

it follows that H has matrix elements only between $|4_n\rangle$ and $|4_n\rangle, |4_n\rangle$ and $|4_{n+1}\rangle$ and $|4_n\rangle$ and $|4_{n+1}\rangle$. Indeed,

$$\langle 4_n | H | 4_n \rangle = Q_{nn}$$

and

$$\langle 4_n | H | 4_{n+1} \rangle (= \langle 4_{n+1} | H | 4_n \rangle^*) = Q_{n+1,n}.$$

Thus the Lanczos algorithm brings H into tridiagonal form.

H can be diagonalized following the procedure described before.

Notice :

- if only the eigenvalues are sought, it is not necessary to keep the vectors $|4_{n-2}\rangle \dots |4_1\rangle$ in memory. Only the current $|4_n\rangle$ and $|4_{n-1}\rangle$ are needed to calculate $|4_{n+1}\rangle$ and only the coefficients a_{nn} and $a_{n-1,n}$ need be stored. This can be particularly useful when dealing with very large, but very sparse matrices.
- the construction of the $|4_n\rangle$ can be truncated, the lowest eigenvectors (i.e. the space spanned by $|4_1\rangle \dots |4_{n_{\max}}\rangle$) can be found, and the procedure can be repeated starting from the lowest eigenvectors. Generally this converges to the true lowest eigenvectors of H . This allows one to find the lowest eigenvalues and eigenvectors of a very large, sparse $N \times N$ matrix H without having to deal with actual $N \times N$ matrices and a characteristic equation of degree N .

Alternative possibilities:

Let the eigenvectors of M be such that

$$\lambda_1 \leq \lambda_2 \dots \leq \lambda_n$$

(or, if complex $\operatorname{Re} \lambda_1 \leq \operatorname{Re} \lambda_2 \dots$).

If only the eigenvector corresponding to the lowest eigenvalue λ_1 , or the eigenvector corresponding to the λ of largest modulus (let this be λ_n) are sought, one can use relaxation techniques.

For instance, the sequence

$$\vec{U}_i \rightarrow \vec{U}_{i+1} = M \vec{U}_i \rightarrow \frac{\vec{U}_{i+1}}{\|\vec{U}_{i+1}\|} \quad \left(\begin{array}{l} \text{this last} \\ \text{step to keep} \\ \text{the vector} \\ \text{normalized} \\ \text{to 1} \end{array} \right)$$

will converge to \vec{U}_n .

Vice versa, if we approximate the evolution

$$\frac{d\vec{U}}{dt} = -M\vec{U}$$

with

$$\vec{U}_i \rightarrow \vec{U}_{i+1} = \vec{U}_i - \epsilon M \vec{U}_i \rightarrow \frac{\vec{U}_{i+1}}{\|\vec{U}_{i+1}\|}$$

and take care that the step ϵ is small enough that

$$\|I - \epsilon \lambda_i\| > \|I - \epsilon \lambda_n\|,$$

then the sequence will converge to the eigenvector corresponding to λ_1 .

Using the orthogonality properties of eigenvectors, these methods can be generalized so as to produce a few of the largest (or smallest) eigenvectors. But they are not well suited to find a large number of eigenvectors.

