

Monday, 01/16/12:

Goals for today:

BLAS = basic linear algebra subroutines

HW trapezoidal rule

HW power method

BLAS = Fortran functions and subroutines

= effort in 1970s to write high-quality, validated, free-of-charge code for basic linear algebra tasks, to be downloaded and incorporated into research and commercial code.

Written by the very best and most famous researchers.

BLAS 1 = level 1 BLAS = one for loop

Assume vectors and matrices have dimension n or $n \times n$, etc.

dot product of vectors

Euclidean norm

vector addition, scaling of a vector

$\alpha x + y$ = scalar α times x plus y = $\alpha x + y$

(\Rightarrow) $\alpha x_i + y_i$ for each i \rightarrow in CPU!

2 operations (αx_i and $+y_i$) in one clock cycle

certain matrix operations that touch only n elements, like

scaling of a row or column

BLAS 2 = level 2 BLAS = 2 (nested) for loops

matrix-vector product

scaling of matrix

outer product = $x * y^T$ = col x times row y^T (\Leftrightarrow) $A_{ij} = x_i y_j$

BLAS 3 = level 3 BLAS = 3 (nested) for loops

matrix-matrix multiplication

LU factorisation \rightarrow LINPACK, LAPACK

packages with higher-level lin. alg.:

LINPACK = lin. alg. like LU, band substitution
EISPACK = eigenvalue algorithms } 1970s, 1980s

1990s there were updated and became LAPACK.

Key: LAPACK is programmed with variable block sizes, depending on size of cache on the actual computer used.

Really use BLAS by linking to a library like ACML

= AMD Core Math Library, sold along with the
Portland Group compiler

libacml.so, header file acml.h \rightarrow See HD and posted code

Documentation: - books from SIAM = Society for Industrial and Applied Math
- webpage www.netlib.org to see actual code and its comments.

At this page, browse for blas, linpack, etc.

We looked at `ddot.f` = double-precision dot product in BLAS

Notice that you can supply increments other than 1 in the call.

Most typical case would have increments of 1, so dot product of vectors x and y would be

$$d = \text{ddot}(n, x, 1, y, 1)$$

Main efficiency for CPUs with a cache lies in loop unrolling

where we pre-load several consecutive components of x and y and act on them simultaneously:

roughly:
$$d = \sum_{i=1:5:n} (x_i y_i + x_{i+1} y_{i+1} + x_{i+2} y_{i+2} + x_{i+3} y_{i+3} + x_{i+4} y_{i+4})$$

How to use increments other than 1 and H/W in power:

$$y = Ax \quad (\Rightarrow) \text{classical thinking } y_i = \sum_{j=1}^n A_{ij} x_j$$

(\Rightarrow) $y_i = \text{dot product of } i\text{-th of } A \text{ with vector } x$

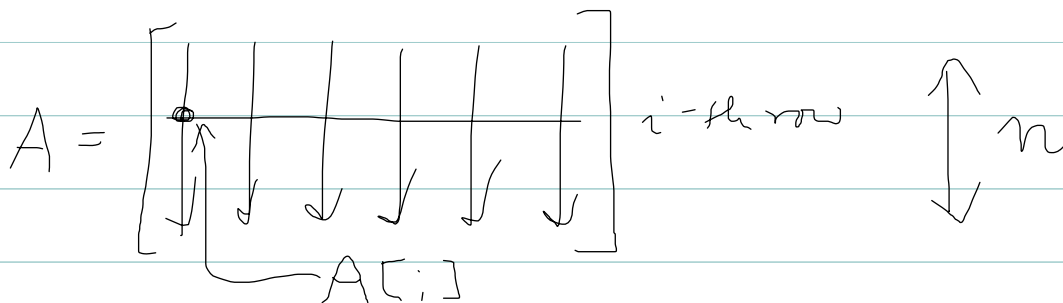
really bad: aux. vectors for the i th: $a_i = A(i,:)$

then $y_i = \text{dot}(a_i, x) = \text{ddot}(n, a_i, 1, x, 1)$

better = no aux. vectors;

for $i = 1:n$

$$y_i = \sum_{j=1}^n A_{ij} x_j = \text{ddot}(n, A[i], n, x, 1)$$



this showed idea. Precise Code

$$d = \text{ddot}(n, \&(A[i]), n, x, 1);$$

If you follow idea above for
 power method / matrix-vector product
 in parallel, you will need one
 parallel_dot call for each row
 of A , that is, n MPI calls
 via MPI_Allreduce in parallel_dot
 That is worse than only 1 or 2
MPI calls?

Find better solution by drawing data:

$$y = Ax = \left[\begin{array}{c|c|c} l_A & l_A & l_A \end{array} \right] \begin{bmatrix} l_x \\ l_x \\ l_x \end{bmatrix}$$

$$= \underbrace{l_A l_x}_{\text{on } P_0} + \underbrace{l_A l_x}_{\text{on } P_1} + \underbrace{l_A l_x}_{\text{on } P_2}$$

There are matrix-vector products on each
 process \Rightarrow we better code for them!

This product:

bad: for $i = 1:n$
for $j = 1:n$

since this accesses A_{ij} out of order and you get cache misses

good: for $j = 1:n$
for $i = 1:n$

so that A 's elements are accessed consecutively.

even better as in BLAS:

idea of block matrices

$$y = Ax = \begin{bmatrix} \boxed{} & \boxed{} & \boxed{} \\ \boxed{} & \boxed{} & \boxed{} \\ \boxed{} & \boxed{} & \boxed{} \end{bmatrix} \begin{bmatrix} \\ \\ \end{bmatrix}$$

(out of time: HW on trapezoidal rule.)