CS601: Software Development for Scientific Computing

Autumn 2023

Week6: Matrix Computations with Sparse Matrices, Tools for debugging and more

LAPACK – Linear Algebra Package

- LAPACK uses BLAS-3 (1989 now)
 - Ex: Obvious way to express Gaussian Elimination (GE) is adding multiples of one row to other rows – BLAS-1
 - How do we reorganize GE to use BLAS-3?
 - Contents of LAPACK (summary)
 - Algorithms that are (nearly) 100% BLAS-3
 - Linear Systems, Least Squares
 - Algorithms that are only ≈50% BLAS-3
 - Eigenproblems, Singular Value Decomposition (SVD)
 - Generalized problems (eg Ax = I Bx)
 - Error bounds for everything
 - Lots of variants depending on A's structure (banded, A=A^T, etc.)
 - How much code? (Release 3.9.0, Nov 2019) (www.netlib.org/lapack)
 - Source: 1982 routines, 827K LOC, Testing: 1210 routines, 545K LOC

Matrix Data and Efficiency

- Sparse Matrices
 - E.g. banded matrices
 - Diagonal
 - Tridiagonal etc.
- Symmetric Matrices

Admit optimizations w.r.t.

- Storage
- Computation

Sparse Matrices - Motivation

 Matrix Multiplication with Upper Triangular Matrices (C=C+AB)

$$\begin{bmatrix} a_{11} & a_{12} & a_{13} \\ 0 & a_{22} & a_{23} \\ 0 & 0 & a_{33} \end{bmatrix} \begin{bmatrix} b_{11} & b_{12} & b_{13} \\ 0 & b_{22} & b_{23} \\ 0 & 0 & b_{33} \end{bmatrix} =$$

$$\begin{bmatrix} A & B & B & B \end{bmatrix}$$

$$\begin{bmatrix} a_{11}b_{11} & a_{11}b_{12} + a_{12} b_{22} & a_{11}b_{13} + a_{12} b_{23} + a_{13} b_{13} \\ 0 & a_{22}b_{22} & a_{22}b_{23} + a_{23} b_{33} \\ 0 & 0 & a_{33}b_{33} \end{bmatrix}$$

A*B

The result, A*B, is also upper triangular.

The non-zero elements appear to be like the result of *inner-product*

Sparse Matrices - Motivation

 C=C+AB when A, B, C are upper triangular, pseudocode: for i=1 to N

- Cost = $\sum_{i=1}^{N} \sum_{j=i}^{N} 2(j-i+1)$ flops (why 2?)
- Using $\Sigma_{i=1}^{N} i \approx \frac{n^2}{2}$ and $\Sigma_{i=1}^{N} i^2 \approx \frac{n^3}{3}$
- $\Sigma_{i=1}^N \Sigma_{j=i}^N 2(j-i+1) \approx \frac{n^3}{3}$, 1/3rd the number of flops required for dense matrix-matrix multiplication

Sparse Matrices

Have lots of zeros (a large fraction)

```
        X
        X
        0
        0
        X
        0
        0
        X

        0
        X
        0
        0
        X
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
        0
```

- Representation
 - Many formats available
 - Compressed Sparse Row (CSR)

```
Implementation:Three arrays:
double *val;
int *ind;
int *rowstart;
```

Sparse Matrices - Example

Using Arrays

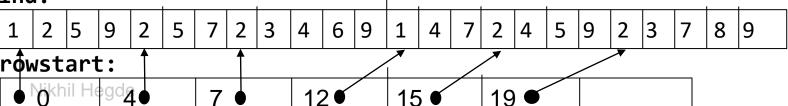
 A_{11} A_{12} A_{12} <t

double *val; //size= NNZ
int *ind; //size=NNZ
int *rowstart; //size=M=Number of rows

val:

						1																		
							l												l	l				
	つ し			a	a	1 1	1	1 1	2	1 ~	1	1	1	1	1	a	1	1	1	1	1	1	1	1
-1-6	a al	a_{12}	dı⊲⊢	d ₄	daa	ldarl	ldaa	เสาร	daal	l d a al	ldac	lana	d 11	l di a al	d₁⊸	a ₅₂	d ₋	ld	la - a	l a ca	$\mathbf{d}_{\mathbf{c}_{\mathbf{a}}}$	ld c ¬	aco	dca
- []	T II	T \	T2	- 19	- //	ı 25	1 - 2 /	1 - 32	- 33	1 - 34	1 - 36	15.39	41	44	-4/	- 52	54	15.55	15.59	15.62	15.63	19.67	15.68	15.69

ind:



Gaxpy with Sparse Matrices: y=y+Ax

Using arrays

```
for i=0 to numRows
  for j=rowstart[i] to rowstart[i+1]-1
  y[i] = y[i] + val[j]*x[ind[j]]
```

- Does the above code reuse y, x, and val ? (we want our code to reuse as much data elements as possible while they are in fast memory):
 - y? Yes. Read and written in close succession.
 - x? Possible. Depends on how data is scattered in val.
 - val? Good spatial locality here. Less likely for a sparse matrix in general.

Nikhil Hegde

Gaxpy with Sparse Matrices: y=y+Ax

Optimization strategies:

```
for i=0 to numRows
  for j=rowstart[i] to rowstart[i+1]-1
  y[i] = y[i] + val[j]*x[ind[j]]
```

- Unroll the j loop // we need to know the number of non-zeros per row
- Eliminate ind[i] and thereby the indirect access to elements of x.
 Indirect access is not good because we cannot predict the pattern of data access in x. //We need to know the column numbers
- Reuse elements of x //The elements of a should be e.g. located closely

These optimizations will not work for y=y+Ax pseudocode in general. When you know the data pattern and metadata info as mentioned above, you can reorder computations (scheduling optimization), reorganize data for better locality.