#### Lecture 10 Sparse Matrices, Iterative Methods

#### David Semeraro

University of Illinois at Urbana-Champaign

September 26, 2013

Latent semantic analysis (LSA) analyzes two-mode data. Looks at relationships between documents and terms.

- natural language processing
- information retrieval
- information filtering
- textual machine learning

Document-term matrix: Document1(D1) = "I love numerical analysis" Document1(D2) = "I do not love numerical analysis, but I love linear algebra."

	I	love	numerical	linear	algebra
D1	1	1	1	0	0
D2	1	2	1	1	1

## An application

	I	love	numerical	linear	algebra
D1	1	1	1	0	0
D2	1	2	1	1	1

One method for weights: Term Count Model

Variation: Term Frequency-Inverse Document Frequency; weight the entries inversely, highlighting infrequent terms

Let *X* be the matrix of occurrences (or the inverse).

$$X = \begin{bmatrix} x_{1,1} & \dots & x_{1,n} \\ \vdots & \ddots & \vdots \\ x_{m,1} & \dots & x_{m,n} \end{bmatrix}$$

Now each row will be a vector relating a term to all documents. Each column will be a vector relating a document to all terms.

$$X = \begin{bmatrix} x_{1,1} & \dots & x_{1,n} \\ \vdots & \ddots & \vdots \\ x_{m,1} & \dots & x_{m,n} \end{bmatrix}$$

- X has many zeros
- a dot product of the rows gives the correlation between terms over the documents
- XX<sup>T</sup> gives a cumulative view of the correlation
- same with X<sup>T</sup>X
- singular value decompositions, eigenvalue analysis, etc give other information

## **Sparse Matrices**

ack: Y. Saad



- Vague definition: matrix with few nonzero entries
- For all practical purposes: an  $m \times n$  matrix is sparse if it has  $O(\min(m, n))$  nonzero entries.
- This means roughly a constant number of nonzero entries per row and column

## **Sparse Matrices**

ack: Y. Saad

- Other definitions use a slow growth of nonzero entries with respect to *n* or *m*.
- Wilkinson's Definition: "..matrices that allow special techniques to take advantage of the large number of zero elements." (J. Wilkinson)"
- A few applications which lead to sparse matrices: Structural Engineering, Computational Fluid Dynamics, Reservoir simulation, Electrical Networks, optimization, data analysis, information retrieval (LSI), circuit simulation, device simulation, ...

#### Sparse Matrices: The Goal

- To perform standard matrix computations economically i.e., without storing the zeros of the matrix.
- For typical Finite Element /Finite difference matrices, number of nonzero elements is O(n).

#### Example

To add two square dense matrices of size *n* requires  $O(n^2)$  operations. To add two sparse matrices *A* and *B* requires O(nnz(A) + nnz(B)) where nnz(X) = number of nonzero elements of a matrix *X*.

#### remark

 $A^{-1}$  is usually dense, but L and U in the LU factorization may be reasonably sparse (if a good technique is used).

Image: Image:

#### Goal

• Principle goal: solve

$$Ax = b$$

where  $A \in \mathbb{R}^{n \times n}$ ,  $x, b \in \mathbb{R}^n$ 

- Assumption: A is very sparse
- General approach: iteratively improve the solution
- Given *x*<sub>0</sub>, ultimate "correction" is

$$x_1 = x_0 + e_0$$

where 
$$e_0 = x - x_0$$
, thus  $Ae_0 = Ax - Ax_0$ ,

or

$$x_1 = x_0 + A^{-1}r_0$$

where  $r_0 = b - Ax_0$ 

- Principle difficulty: how do we "approximate"  $A^{-1}r$  or reformulate the iteration?
- One simple idea:

$$x_1 = x_0 + \alpha r_0$$

- operation is inexpensive if r<sub>0</sub> is inexpensive
- requires very fast sparse mat-vec (matrix-vector multiply) Ax<sub>0</sub>

- So how do we store *A*?
- Fast mat-vec is certainly important; also ask
  - what type of access (rows, cols, diag, etc)?
  - dynamic allocation?
  - transpose needed?
  - inherent structure?
- Unlike dense methods, not a lot of standards for iterative
  - dense BLAS have been long accepted
  - sparse BLAS still iterating
- Even data structures for dense storage not as obvious
- Sparse operations have low operation/memory reference ratio

Matrix Market attempts to classify the sparse matrix. First Qualification (type of values and number of values):

identifier	description
Real	All entries are float
Complex	All entries are a pair of float
Integer	All entries are int
Pattern	Matrix is a pattern. Actual entries are omitted
Parallel	Parallel structure is identified

Image: A matrix and a matrix

Second Qualification (interpreting values):

description
A has no symmetry, no symmetry is utilized,
or $A$ is not square
$a_{ij} = a_{ji}$ ; only entries on the diagonal
and below(or above) are stored.
$a_{ii} = -a_{ii}$ ; only entries below (or above)
the diagonal (= 0) are stored.
$a_{ii} = \bar{a}_{ii}$ ; only entries on the diagonal
and below (or above) are stored.

see "The Matrix Market Exchange Formats: Initial Design" by Boisvert, Pozo, Remington

イロト イポト イヨト イヨト

- DNS Dense
- BND Linpack Banded
- COO Coordinate
- **CSR** Compressed Sparse Row
- **CSC** Compressed Sparse Column
- MSR Modified CSR
- LIL Linked List

- ELL Ellpack-Itpack
- DIA Diagonal
- BSR Block Sparse Row
- **SSK** Symmetric Skyline
- **BSR** Nonsymmetric Skyline
- JAD Jagged Diagonal

note: CSR = CRS, CCS = CSC, SSK = SKS in some references

$$A = \begin{bmatrix} 1.0 & 2.0 & 3.0 \\ 4.0 & 5.0 & 6.0 \\ 7.0 & 8.0 & 9.0 \end{bmatrix}$$
$$AA = \begin{bmatrix} 3 & 3 & 1.0 & 2.0 & 3.0 & 4.0 & 5.0 & 6.0 & 7.0 & 8.0 & 9.0 \end{bmatrix}$$

- simple
- row-wise
- easy blocked formats

≣ ▶

<ロト < 回ト < 回ト

Ð

I

590

$$A = \begin{bmatrix} 1 & 0 & 0 & 2 & 0 \\ 3 & 4 & 0 & 5 & 0 \\ 6 & 0 & 7 & 8 & 9 \\ 0 & 0 & 10 & 11 & 0 \\ 0 & 0 & 0 & 0 & 12 \end{bmatrix}$$
$$AA = \begin{bmatrix} 12.0 & 9.0 & 7.0 & 5.0 & 1.0 & 2.0 & 11.0 & 3.0 & 6.0 & 4.0 & 8.0 & 10.0 \\ JR = \begin{bmatrix} 5 & 3 & 3 & 2 & 1 & 1 & 4 & 2 & 3 & 2 & 3 & 4 \\ JC = \begin{bmatrix} 5 & 5 & 3 & 4 & 1 & 4 & 4 & 1 & 1 & 2 & 4 & 3 \end{bmatrix}$$

• simple, often used for entry

Question: Do you need this much storage?



$$A = \begin{bmatrix} 1 & 0 & 0 & 2 & 0 \\ 3 & 4 & 0 & 5 & 0 \\ 6 & 0 & 7 & 8 & 9 \\ 0 & 0 & 10 & 11 & 0 \\ 0 & 0 & 0 & 0 & 12 \end{bmatrix}$$

$$AA = \begin{bmatrix} 1 & 0 & 2.0 & 3.0 & 4.0 & 5.0 & 6.0 & 7.0 & 8.0 & 9.0 & 10.0 & 11.0 & 12.0 \\ JA = \begin{bmatrix} 1 & 4 & 1 & 2 & 4 & 1 & 3 & 4 & 5 & 3 & 4 & 5 \\ IA = \begin{bmatrix} 1 & 3 & 6 & 10 & 12 & 13 \end{bmatrix}$$

- Length of AA and JA is nnz; length of IA is n + 1
- *IA*(*j*) gives the index (offset) to the beginning of row *j* in *AA* and *JA* (one origin due to Fortran)
- no structure, fast row access, slow column access (why?)
- related: CSC, MSR



AA

$$A = \begin{bmatrix} 1 & 0 & 0 & 2 & 0 \\ 3 & 4 & 0 & 5 & 0 \\ 6 & 0 & 7 & 8 & 9 \\ 0 & 0 & 10 & 11 & 0 \\ 0 & 0 & 0 & 0 & 12 \end{bmatrix}$$
$$AA = \begin{bmatrix} 1.0 & 4.0 & 7.0 & 11.0 & 12.0 & * & 2.0 & 3.0 & 5.0 & 6.0 & 8.0 & 9.0 & 10.0 \end{bmatrix}$$
$$JA = \begin{bmatrix} 7 & 8 & 10 & 13 & 14 & 14 & 4 & 1 & 4 & 1 & 4 & 5 & 3 \end{bmatrix}$$

- places importance on diagonal (often nonzero and accessed frequently)
- first n entries are the diag
- n+1 is empty
- rest of AA are the nondiagonal entries
- first n+1 entries in *IA* give the index (offset) of the beginning of the row (the IA of CSR is in this IA)
- rest of IA are the columns indices

$$A = \begin{bmatrix} 1 & 0 & 2 & 0 & 0 \\ 3 & 4 & 0 & 5 & 0 \\ 0 & 6 & 7 & 0 & 8 \\ 0 & 0 & 9 & 10 & 0 \\ 0 & 0 & 0 & 11 & 12 \end{bmatrix} \quad DIAG = \begin{bmatrix} * & 1.0 & 2.0 \\ 3.0 & 4.0 & 5.0 \\ 6.0 & 7.0 & 8.0 \\ 9.0 & 10.0 & * \\ 11.0 & 12.0 & * \end{bmatrix} \quad IO$$

$$IOFF = \begin{bmatrix} -1 & 0 & 2 \end{bmatrix}$$

- need to know the offset structure
- some entries will always be empty



- Form columns from first non-zero in each row, repeat.
- used more on vector machines (what? why?)
- assumes low number of *nnz* per row (=number of columns in *COEFF* and *JCOEFF*)

#### JDS

- like DIA (and CDS),but more space-efficient
- costs more to gather and scatter

Take

$$A = \begin{bmatrix} 10 & -3 & 0 & 1 & 0 & 0 \\ 0 & 9 & 6 & 0 & -2 & 0 \\ 3 & 0 & 8 & 7 & 0 & 0 \\ 0 & 6 & 0 & 7 & 5 & 4 \\ 0 & 0 & 0 & 0 & 9 & 13 \\ 0 & 0 & 0 & 0 & 5 & -1 \end{bmatrix}$$

And shift:

$$A \rightarrow \begin{bmatrix} 10 & -3 & 1 \\ 9 & 6 & -2 \\ 3 & 8 & 7 \\ 6 & 7 & 5 & 4 \\ 9 & 13 \\ 5 & -1 \end{bmatrix}$$

$$A \rightarrow \begin{bmatrix} 10 & -3 & 1 \\ 9 & 6 & -2 \\ 3 & 8 & 7 \\ 6 & 7 & 5 & 4 \\ 9 & 13 \\ 5 & -1 \end{bmatrix}$$

Now store the columns and column indices:

$$VAL = \begin{bmatrix} 10 & 9 & 3 & 6 & 9 & 5 \\ -3 & 6 & 8 & 7 & 13 & -1 \\ 1 & -2 & 7 & 5 & 0 & 0 \\ 0 & 0 & 0 & 4 & 0 & 0 \end{bmatrix} \quad COL = \begin{bmatrix} 1 & 2 & 1 & 2 & 5 & 5 \\ 2 & 3 & 3 & 4 & 6 & 6 \\ 4 & 5 & 4 & 5 & 0 & 0 \\ 0 & 0 & 0 & 6 & 0 & 0 \end{bmatrix}$$

イロト イロト イヨト イヨト

- now reorder in terms of largest to smallest nnz in each row
- in JDS, the number of jagged diagonals is the number of *nnz* in the first row of *VAL*. This is the max *nnz* in any row.
- PERM: permutation array to reorder rows
- JDIAG: jagged diags in order
- COL: column indices
- JDPTR: points to the beginning of each diagonal
- advantage for mat-mat (see "Krylov subspace methods on supercomputers" by Saad)
- this is actually ITPACK or Purdue storage

$$JDIAG = \begin{bmatrix} 6 & 9 & 3 & 10 & 9 & 5; & 7 & 6 & 8 & -3 & 13 & -1; & 5 & -2 & 7 & 1; & 4 \\ COL = \begin{bmatrix} 2 & 2 & 1 & 1 & 5 & 5; & 4 & 3 & 3 & 2 & 6 & 6; & 5 & 5 & 4 & 4; & 6 \\ PERM = \begin{bmatrix} 4 & 2 & 3 & 1 & 5 & 6 \end{bmatrix}$$
  
$$JDPTR = \begin{bmatrix} 1 & 7 & 13 & 17 \end{bmatrix}$$

## Blocked

$$A = \begin{bmatrix} 1.0 & 2.0 & 0.0 & 0.0 & 3.0 & 4.0 \\ 5.0 & 6.0 & 0.0 & 0.0 & 7.0 & 8.0 \\ 0.0 & 0.0 & 9.0 & 10.0 & 11.0 & 12.0 \\ 0.0 & 0.0 & 13.0 & 14.0 & 15.0 & 16.0 \\ 17.0 & 18.0 & 0.0 & 0.0 & 20.0 & 21.0 \\ 22.0 & 23.0 & 0.0 & 0.0 & 24.0 & 25.0 \end{bmatrix}$$
$$AA = \begin{bmatrix} 1.0 & 3.0 & 9.0 & 11.0 & 17.0 & 20.0 \\ 5.0 & 7.0 & 15.0 & 13.0 & 22.0 & 24.0 \\ 2.0 & 4.0 & 10.0 & 12.0 & 18.0 & 21.0 \\ 6.0 & 8.0 & 14.0 & 16.0 & 23.0 & 25.0 \end{bmatrix}$$
$$JA = \begin{bmatrix} 1 & 5 & 3 & 5 & 5 & 1 & 5 \\ IA = \begin{bmatrix} 1 & 5 & 3 & 5 & 5 & 1 & 5 \end{bmatrix}$$
$$IA = \begin{bmatrix} 1 & 3 & 5 & 7 \end{bmatrix}$$

David Semeraro (NCSA)

September 26, 2013 23 / 38

1

< □ > < □ > < □ > < □ > < □ >

**]** האר

- each column of AA is a 2 × 2 block
- JA(k) = column index of (1, 1) entries of the *kth* block
- declared as AA(2, 2, 6)
- blocks arise in many apps
- variant: variable block size

#### Blocked

Also row-wise

$$AA = \begin{bmatrix} 1.0 & 5.0 & 2.0 & 6.0 \\ 3.0 & 7.0 & 4.0 & 8.0 \\ 9.0 & 15.0 & 10.0 & 14.0 \\ 11.0 & 13.0 & 12.0 & 16.0 \\ 17.0 & 22.0 & 18.0 & 23.0 \\ 20.0 & 24.0 & 21.0 & 25.0 \end{bmatrix}$$

$$JA = \begin{bmatrix} 1 & 5 & 3 & 5 & 5 & 1 & 5 \end{bmatrix}$$
$$IA = \begin{bmatrix} 1 & 3 & 5 & 7 \end{bmatrix}$$

- each row of AA is a 2  $\times$  2 block (can be a drawback)
- *JA*, *IA* same, *AA*(6, 2, 2)
- if elements of blocks are accessed at the same time: rows are better (C)
- if elements of similar positions in different blocks are accessed at the same time: columns are better (C)

## SSK, NSK

- for "skyline" matrices (variable band, see "Direct methods for sparse matrices" by Duff, Erisman, Reid)
- can be used for diagonal block matrices
- skyline structure is preserved in basic GE
- for symmetric: Place all the rows (in order) into VAL
- IA points to the beginning of each row
- JA implicit
- for nonsymmetric: store L in SSK format, then U in column-wise SSK see

"SPARSEKIT" by Saad



- Similar to CSR, but rather than a flat AA vector, each row is a linked list of elements
- first element of each row is accessed by ROOT
- each element in AA has a corresponding NEXT entry
- -1 indicates the end of a row
- column lookup take O(nnz); one semi-costly fix: store a columnwise index in the same way as rows.
- very good element insertion time, but more memory

$$A = \begin{bmatrix} 7 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 2 & 0 & 0 & 0 \\ 0 & 2 & 0 & 2 & 0 & 0 \\ 0 & 0 & 0 & 0 & 5 & 0 \\ 0 & 0 & 0 & 0 & 6 & 4 \end{bmatrix}$$

- CSR
- CSC
- COO
- LIL

1

< □ > < □ > < □ > < □ > < □ > .

I

500

## Example

A =	7 0 0 0 0	0 0 1 2 2 0 0 0 0 0	0 0 2 0 0	$\begin{array}{ccc} 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 5 & 0 \\ 6 & 4 \end{array}$	-	<i>i</i> 1 2 3 4 5 6 7 8	<i>IA</i> 2 3 4 2 5 1 5 3	JA 2 4 5 3 6 1 5 2	AA 1 2 5 2 4 7 6 2	000	
i	IA 1	JA 1	<i>AA</i>	_	<i>i</i>	IA	JA	AA 2	NE	<u>XT</u>	
2	2	2	1		2	ד 2	2	2	Ē	5	
3	4	3	2		3	2	2	1	-	) I	
4	6	2	2	CSR	4	8	1	7	-	1	LIL
5	7	4	2		5	6	4	2	-	1	
6	9	5	5		6	-	5	6	7	7	
7	-	5	6		7	-	6	4	-	1	1
8	-	6	4		8	-	_ 5,	چ <sub>ھ</sub> پ	< ≣ ≻ <	<b>1</b> ≞ →	<b>⊥</b> ה∢רי ≣

#### Sparse Matrix-Vector Multiply

$$z = Ax, A_{m \times n}, x_{n \times 1}, z_{m \times 1}$$
  
input A, x  
$$z = 0$$
  
for  $i = 1$  to m  
for  $col = A(i, :)$   
$$z(i) = z(i) + A(i, col)x(col)$$
  
end  
r end

- difference between CSR and LIL is computing line 4
- CSR: rows are contiguous...(next slide)
- LIL: follow rows through linked list

# Sparse Matrix-Vector Multiply

$$z = Ax, A_{m \times n}, x_{n \times 1}, z_{m \times 1}$$

- 1 DO I=1, m
  2 Z(I)=0
  3 K1 = IA(I)
  4 K2 = IA(I+1)-1
  5 DO J=K1, K2
  6 z(I) = z(I) + A(J)\*x(JA(J))
  7 ENDDO
  8 ENDDO
  - O(nnz)
  - marches down the rows
  - very cheap

イロト イヨト イヨト イ

#### Sparse Matrix-Vector Multiply CSR and Data Streams

```
z = Ax, A_{m \times n}, x_{n \times 1}, z_{m \times 1}
_{1} DO I=1, M/2
z = z(I) = 0
 z(I+M/2) = 0
3
 K1 = IA(I)
4
 K2 = IA(I+1) - 1
5
 K3 = IA(I+M/2)
6
   K4 = IA(I+M/2+1) - 1
7
    DO J=0, MIN(K2-K1, K4-K3)
8
       z(I) = z(I) + A(K1+J)*x(JA(K1+J))
9
       z(I+M/2) = z(I+M/2) + A(K3+J)*x(JA(K3+J))
10
    ENDDO
11
   ! ... finish up
12
13 ENDDO
```

#### IA structure allows two data streams

イロト イポト イヨト イヨト

#### Sparse Matrix-Matrix Multiply

• ways to optimize ("SMPP", Douglas, Bank)

```
Z = AB, A_{m \times n}, B_{n \times p}, z_{m \times p}
for i = 1 to m
for j = 1 to n
Z(i,j) = dot(A(i,:), B(:,j))
end
fend
for Z
```

- obvious problem: column selection of B is expensive for CSR
- not-so-obvious problem: Z is sparse(!!), but the algorithm doesn't account for this.

・ロト ・同ト ・ヨト ・ヨ

## Sparse Matrix-Matrix Multiply

$$Z = AB, A_{m \times n}, B_{n \times p}, z_{m \times p}$$

$$\stackrel{1}{=} Z=0$$

$$\stackrel{2}{=} \text{for } i = 1 \text{ to } m$$

$$\stackrel{3}{=} \text{for } colA = A(i, :)$$

$$\stackrel{4}{=} for \ colB = A(colA, :)$$

$$\stackrel{5}{=} Z(i, colB) + = A(i, colA) \cdot B(colA, colB)$$

$$\stackrel{6}{=} \text{end}$$

$$\stackrel{7}{=} \text{end}$$

$$\stackrel{8}{=} \text{end}$$

$$\stackrel{9}{=} \text{return } Z$$

- only marches down rows
- only computes nonzero entries in Z (aside from fortuitous subtractions)
- line 5 will do and insert into Z. Two options:
  - precompute sparsity of Z in CSR
  - use LIL for Z





1

**]** האר

#### Some Python

From COO to CSC:

```
1 from scipy import sparse
2 from numpy import array
3 import pprint
4 IA=array([1,2,3,1,4,0,4,2])
5 JA=array([1,3,4,2,5,0,4,1])
6 V=array([1,2,5,2,4,7,6,2])
7
8 A=sparse.coo_matrix((V,(IA,JA)),shape=(5,6)).tocsr()
```

Nonzeros:

print(A.nnz)

To full and view:

- 1 B=A.todense()
- 2 pprint.pprint(B)

Э

< □ > < □ > < □ > < □ >

#### Simple Matrix Iterations

Solve

$$Ax = b$$

- Assumption: A is very sparse
- Let A = N + M, then

$$Ax = b$$
  
(N+M)x = b  
Nx = b-Mx

Make this into an iteration:

$$Nx_k = b - Mx_{k-1}$$
  
 $x_k = N^{-1}(b - Mx_{k-1})$ 

- Careful choice of N and M can give effective methods
- More powerful iterative methods exist