



**BIG DATA**

**MapReduce Algorithms**

**– Matrix Multiplication**

---

**K V Subramaniam**

Computer Science and Engineering

# Matrices and Vectors - introduction

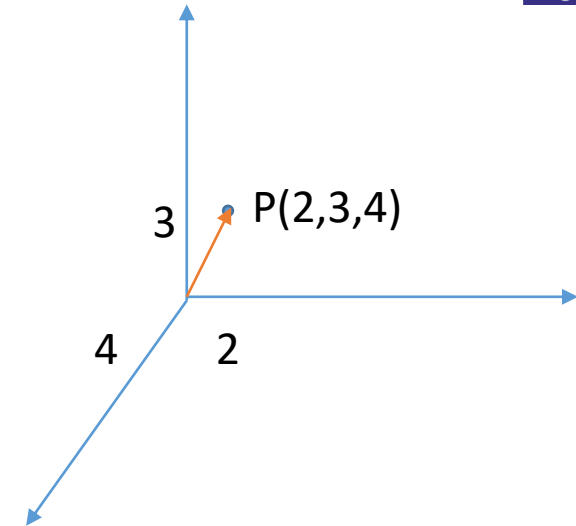
- Matrix Multiplication algorithms
  - Fundamental to many computations, including Page Rank
- Source
  - Leskovec, Jure, Anand Rajaraman, and Jeffrey David Ullman. *Mining of massive datasets*. Cambridge University Press, 2014.
  - <http://infolab.stanford.edu/~ullman/mmds/book.pdf>
  - 4.3.2 of T1

# BIG DATA

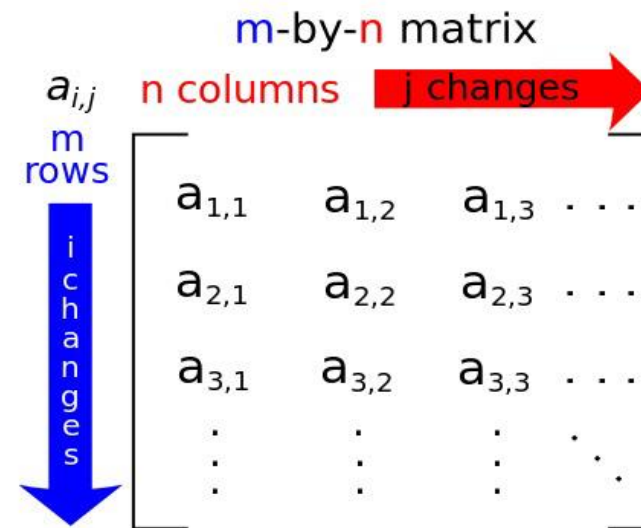
## Background: Vectors and Matrices



- Vectors
  - Can be defined as an ordered list of numbers
  - Visualization
    - An arrow where the direction of the vector is given by the relative size of the components
- Some Common Operations
  - Addition:  $\mathbf{v} + \mathbf{w}$ 
    - Add components
  - Scalar multiplication  $a\mathbf{v}$ 
    - Multiply each component by constant



- Rectangular array of numbers.
- The numbers are called the elements of the matrix.
- An  $m \times n$  matrix has  $m$  rows and  $n$  columns
  - Can be considered as a collection of
    - $m$  row vectors
    - $n$  column vectors
  - An  $n \times n$  matrix is called a square matrix.
- Vector can be considered as a
  - $1 \times n$  matrix (row matrix)
  - $n \times 1$  matrix (column matrix)



Each element of a matrix is often denoted by a variable with two **subscripts**. For example,  $a_{2,1}$  represents the element at the second row and first column of a matrix **A**.

- Multiply each row vector of **A** by the corresponding elements of **x** and sum
- Multiplying a  $m \times n$  matrix by an  $n$  element vector gives an  $m$  element vector

$$A\mathbf{x} = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$
$$= \begin{bmatrix} a_{11}x_1 + a_{12}x_2 + \cdots + a_{1n}x_n \\ a_{21}x_1 + a_{22}x_2 + \cdots + a_{2n}x_n \\ \vdots \\ a_{m1}x_1 + a_{m2}x_2 + \cdots + a_{mn}x_n \end{bmatrix}.$$

$$A = \begin{bmatrix} 1 & 0 & -1 \\ 3 & 1 & 2 \end{bmatrix}, \quad \mathbf{x} = (x, y, z)$$

$$A\mathbf{x} = \begin{bmatrix} 1 & 0 & -1 \\ 3 & 1 & 2 \end{bmatrix} \begin{bmatrix} x \\ y \\ z \end{bmatrix} = \begin{bmatrix} x - z \\ 3x + y + 2z \end{bmatrix}$$
$$= (x - z, 3x + y + 2z).$$

- Typically, matrices are stored as multi-dimensional arrays in programs
- `int A[10][10]`  
allocates 100 integers and is accessed as a 10x10 matrix

Space required to store the matrix – =  
 $10 \times 10 \times \text{sizeof}(\text{int}) = 100 \times \text{sizeof}(\text{int}) = 100 \times 4 = 400$  bytes.

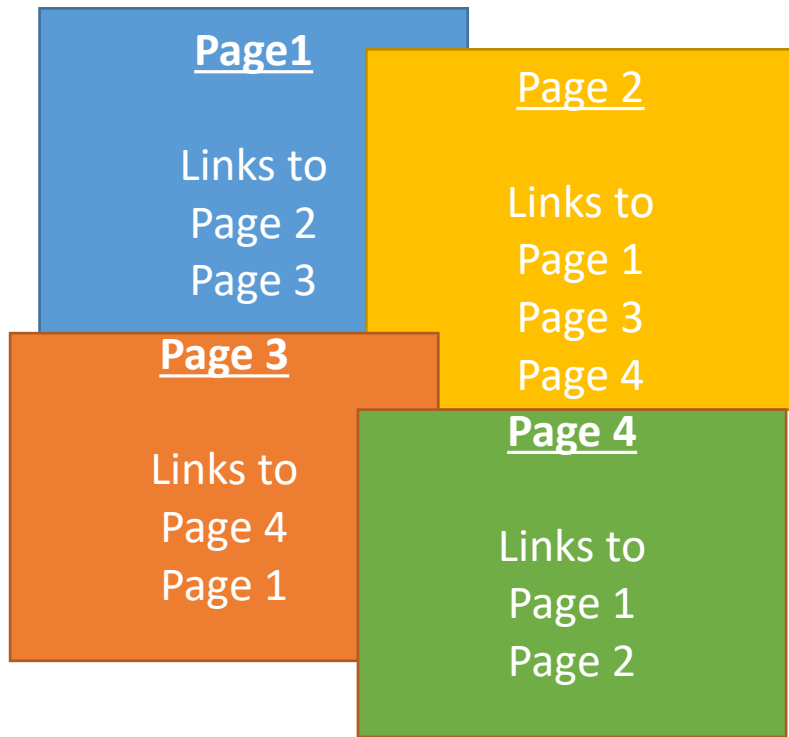
In general, we need  $n^2$  integers to store an ***nxn*** matrix

## Matrix Representation of WWW



# BIG DATA

## How do you represent the pages in WWW?

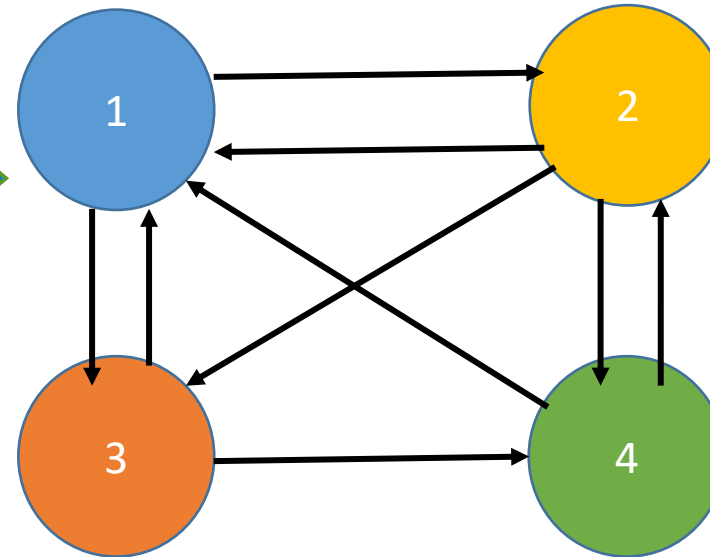
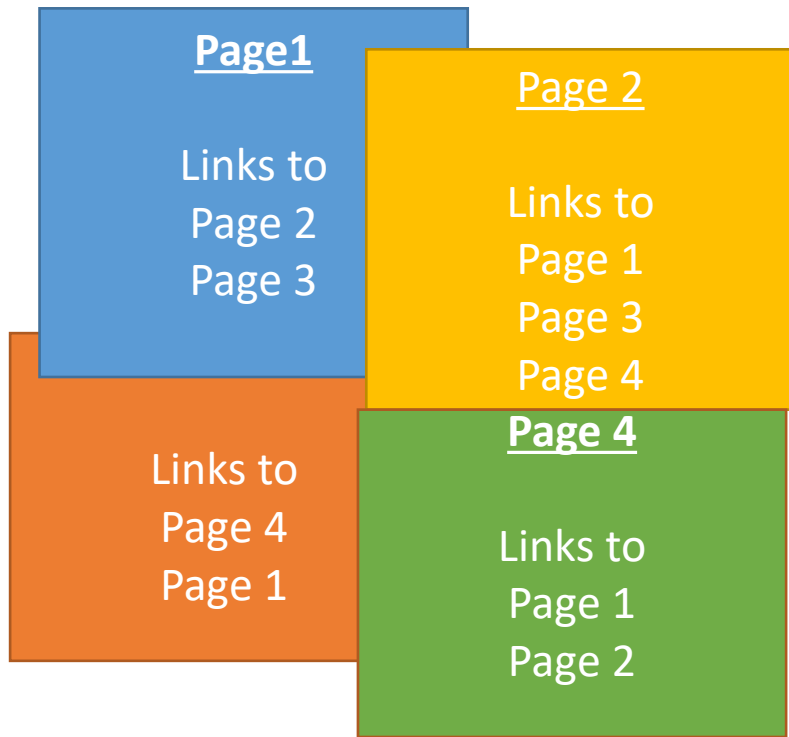


Consider a sample of the internet that contains 4 pages

- How should we represent this?

# BIG DATA

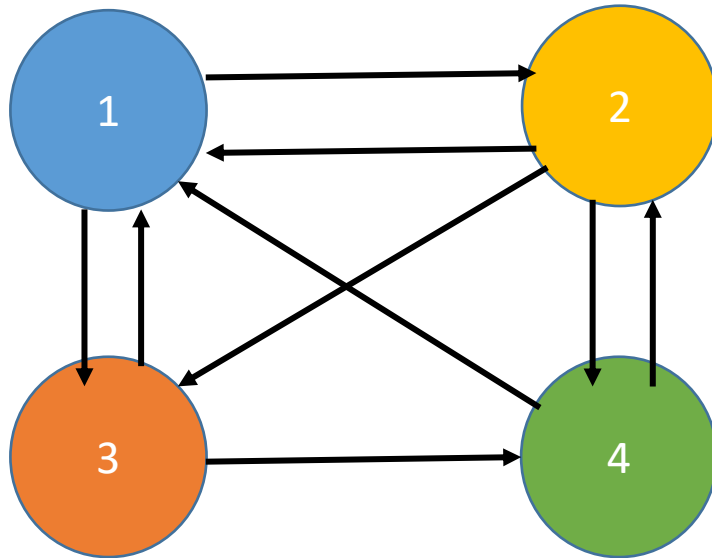
## Modelling the WWW as a directed graph



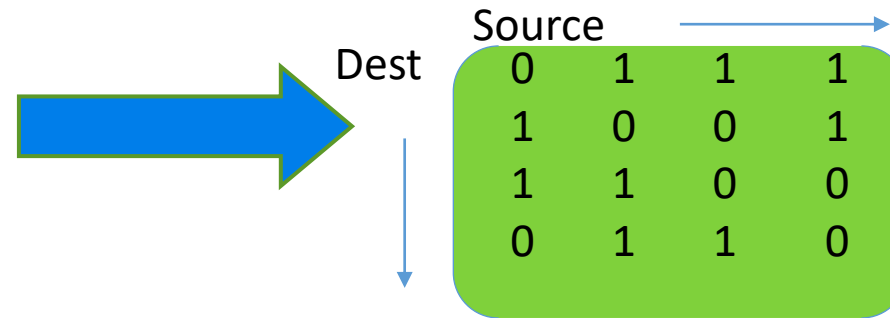
Represented as a directed graph

Pages in the WWW

## Representing the graph as an Adjacency Matrix

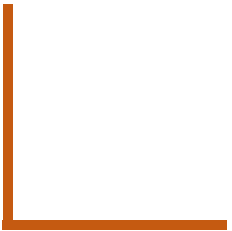


Directed Graph



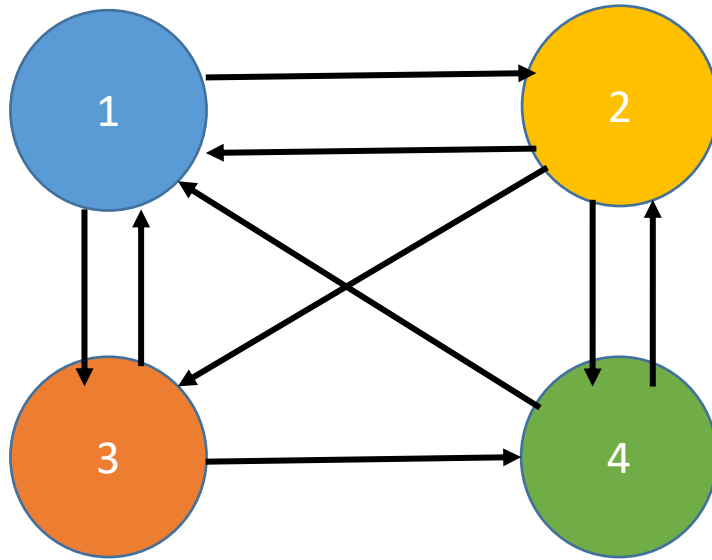
- This is fine for a small graph – 4 pages.
- But internet is large – billions of pages.
- How much storage will we require?

## Large scale matrix representations

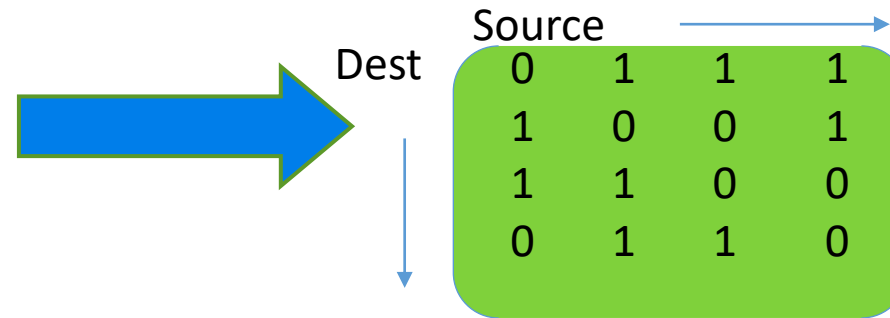


# BIG DATA

## Representing the graph as an Adjacency Matrix



Directed Graph



- Internet is large – billions of pages.
- Note – most of the entries will be 0
- Store as a sparse matrix..

# BIG DATA

## Sparse Matrix representation

---



- In Big Data, we deal with large matrices
  - e.g, n will be the order of  $10^{10}$  if n is number of web pages
- And it will be a sparse matrix
- Won't fit in the memory (DRAM)
- Have to store it in HDFS

- Store only non-zero elements as a separate record in CSV format
- For each element store
  - *<row\_number, column\_number, value>*
  - As the format
- As many entries as there are links
- Exercise –Store the graph given on the right into a HDFS CSV file

0	1	1	1
1	0	0	1
1	1	0	0
0	1	1	0

# BIG DATA

## Solution

---

1, 2, 1  
1, 3, 1  
1, 4, 1  
2, 1, 1  
2, 4, 1  
3, 1, 1  
3, 2, 1  
4, 2, 1  
4, 3, 1

0	1	1	1
1	0	0	1
1	1	0	0
0	1	1	0

As an exercise, try saving this in a file and loading it onto HDFS that you have installed.



# Matrix Vector Multiplication

- To multiply an  $n \times n$  matrix  $M$  with an  $n$ -element vector  $v$ , compute

$$x_i = \sum_{j=1}^n m_{ij} v_j$$

$$A\mathbf{x} = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n \\ a_{21}x_1 + a_{22}x_2 + \dots + a_{2n}x_n \\ \vdots \\ a_{m1}x_1 + a_{m2}x_2 + \dots + a_{mn}x_n \end{bmatrix}$$

- Let us assume that the vector  $\mathbf{v}$  fits into memory
- Vector  $\mathbf{v}$  is shared by all the mappers
- $M_{ij}$  is stored as a CSV file on HDFS and is distributed across multiple nodes

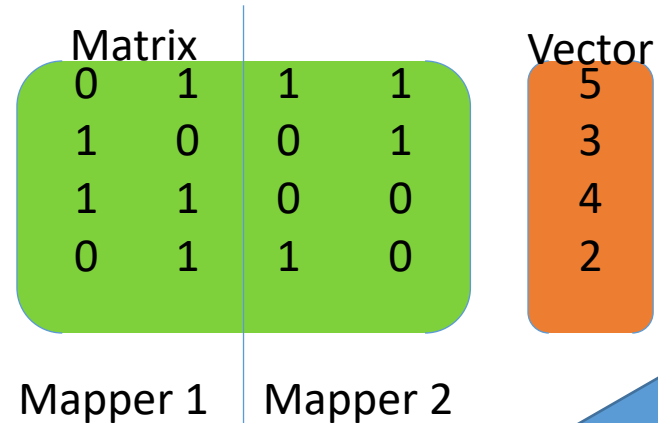
$$x_i = \sum_{j=1}^n m_{ij} v_j$$

$$A\mathbf{x} = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n \\ a_{21}x_1 + a_{22}x_2 + \dots + a_{2n}x_n \\ \vdots \\ a_{m1}x_1 + a_{m2}x_2 + \dots + a_{mn}x_n \end{bmatrix}$$

← 2

$$x_i = \sum_{j=1}^n m_{ij} v_j$$

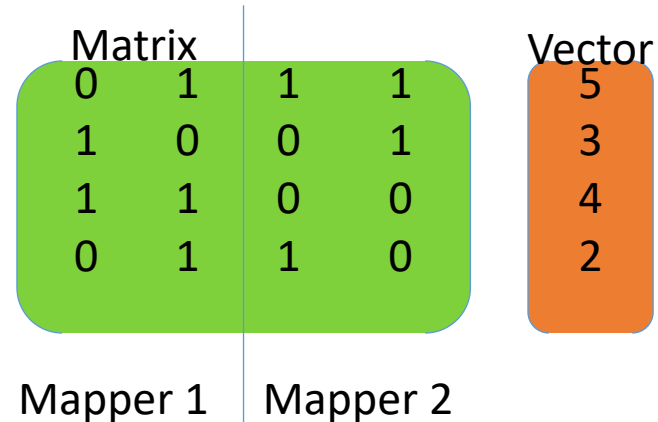
- **map:**
  - Computes the partial product
  - Uses the key as  $i \rightarrow$  the index into the target vector
  - output  $(i, m_{ij} v_j)$
- **reduce:**
  - Sums all the partial products.



Note that the key is the index into the vector where this value will contribute

Key	Value
1	$1*3 = 3$
2	$1*5 = 5$
3	$1*5 = 5$
3	$1*3 = 3$
4	$1*3 = 3$

Key	Value
1	$1*4 = 4$
1	$1*2 = 2$
2	$1*2 = 2$
4	$1*4 = 4$



Reducer Input

Key	Intermediate Value List
1	2,3,4
2	2, 5
3	3, 5
4	3, 4

Reducer output

Key	Value
1	9
1	7
2	8
4	7

Matrix

0	0	3	8
0	9	5	0
0	10	0	0
5	0	0	1

Vector

1  
2  
3  
4

- Assuming rows 1 and 3 are in datanode1 and 2 and 4 are in datanode 2, perform a matrix multiplication using Map Reduce
- Show inputs/outputs of mappers/reducers as discussed in previous slides

## Matrix Vector Multiplication - extensions



- Case 2:  $\mathbf{v}$  doesn't fit into main memory.
- Partition  $M$  and  $\mathbf{v}$  into stripes.

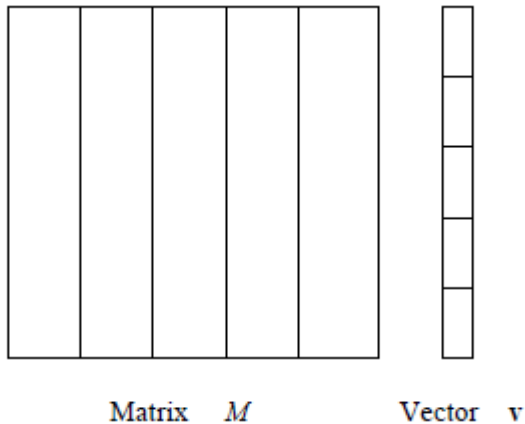


Figure 2.4: Division of a matrix and vector into five stripes

- The same MapReduce algorithm can be used.



# THANK YOU

---

**K V Subramaniam**

Dept. of Computer Science and Engineering

[subramaniamkv@pes.edu](mailto:subramaniamkv@pes.edu),

[ushadevibg@pes.edu](mailto:ushadevibg@pes.edu)