

POLITECNICO
MILANO 1863

Recommender Systems

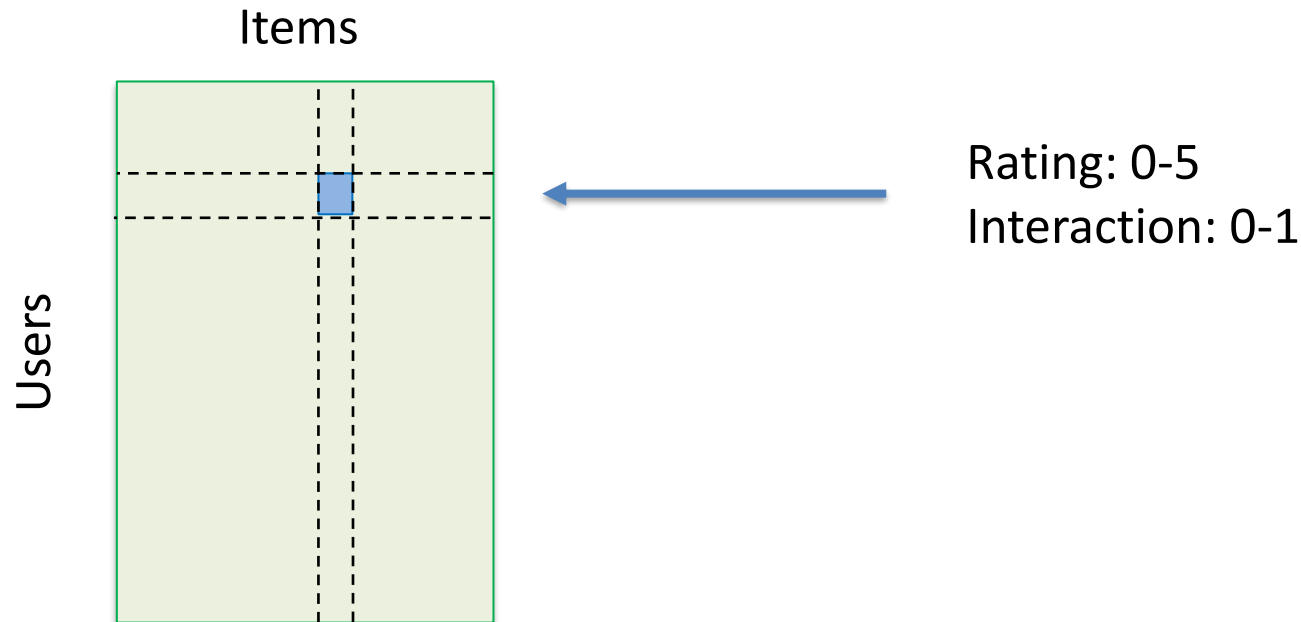
Sparse matrices with Scipy

Maurizio Ferrari Dacrema

A.Y. 2021/22

User Rating Matrix (URM)

In Recommender Systems we often rely on user-item interaction data, which tells us the items any user interacted with.



User Rating Matrix (URM)

The URM has a number of cells equal to $|items| \cdot |users|$

Consider a relatively small case:

- $100'000 = 10^5$ users and items
- 10^{10} cells or $10^{10} \cdot 16 \text{ bit} = 80 \text{ GB}$ to store integer interaction data
- Most of the data will be zero. The number of interactions that occurred is typically in the range of $10^{-2} / 10^{-5}$

Need better solution!

Underlying idea:

- Assume all cells have a default value (zero)
- Store only the cells that have a different value

Challenges:

- We need fast access
- We need fast linear algebra operations, e.g., multiplication

The COOrdinate format stores values as triples (row, col, data)

$$\begin{bmatrix} 5 & & \\ & 4 & \\ 2 & & 3 \end{bmatrix}$$

In scipy:

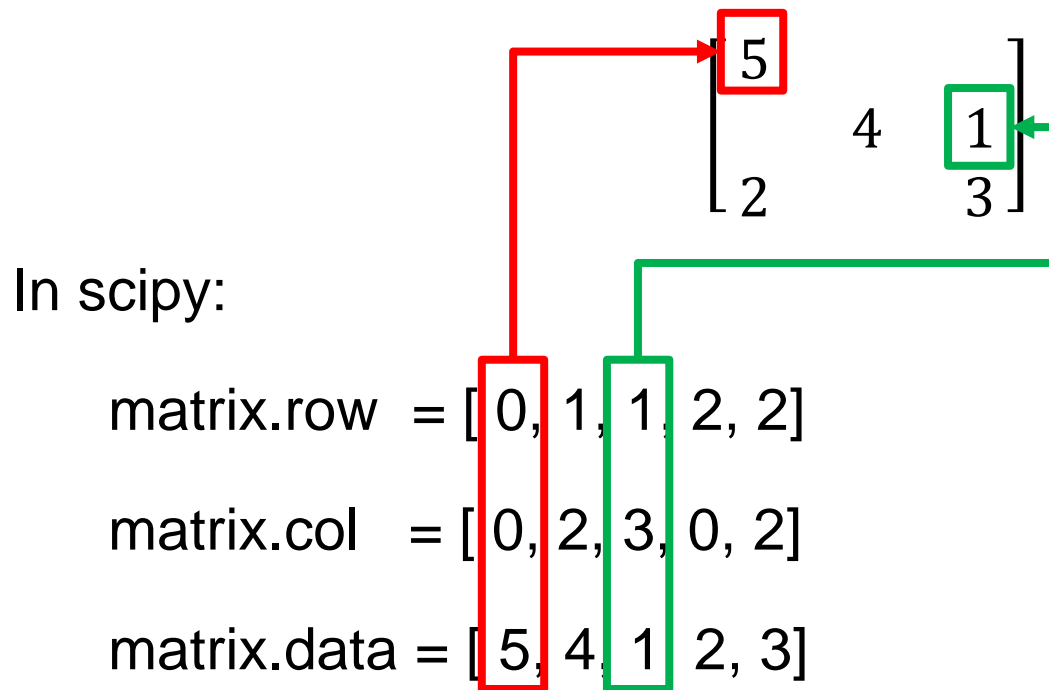
```
matrix.row = [ 0, 1, 1, 2, 2]
```

```
matrix.col = [ 0, 2, 3, 0, 2]
```

```
matrix.data = [ 5, 4, 1, 2, 3]
```

COOrdinate Format

The COOrdinate format stores values as triples (row, col, data)



COO format is easy to read

row: 0,0,0,0,0,0,0,0,1,1,1,1,1,1,2,2,2,3,3,3 ...
col: 1,5,0,2,5,7,0,3,2,0,8,6,0,1,1,2,0,6,5,3 ...

COO format is easy to read but not usable for matrix operations

```
row: 0,0,0,0,0,0,0,0,1,1,1,1,1,1,2,2,2,3,3,3 ...  
col: 1,5,0,2,5,7,0,3,2,0,8,6,0,1,1,2,0,6,5,3 ...
```

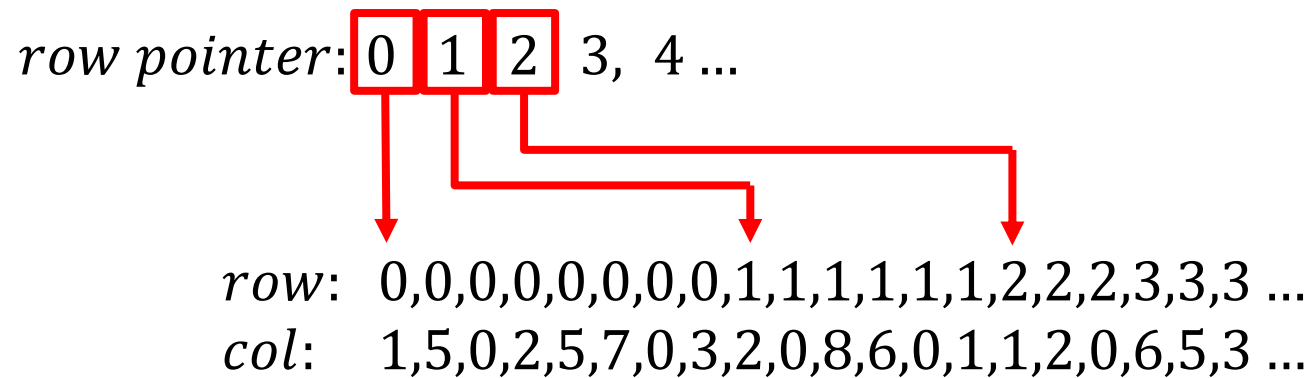
How to find the elements in column 5? Loop all the “col” array

How to find row 2569? Loop the “row” array until you find it

Impossible to do fast access and matrix operations

From COO to Compressed formats

What if we store separately a pointer to where a row begins?



Compressed Sparse Row (CSR) format

CSR format is developed to provide fast row access

The idea:

- Store only the (col, data) explicitly but not the row
- Use a “indptr” data structure to know where the row begins in the (col, data) structure

Data for row i will be between positions $indptr[i]:indptr[i + 1]$

Compressed Sparse Row (CSR) format

$$\begin{bmatrix} 5 & & \\ & 4 & 1 \\ 2 & & 3 \end{bmatrix}$$

In scipy:

```
matrix.indptr = [ 0, 1, 3, 5 ]
```

```
matrix.indices = [ 0, 1, 2, 0, 2 ]
```

```
matrix.data = [ 5, 4, 1, 2, 3 ]
```

Compressed Sparse Row (CSR) format

$$\begin{bmatrix} 5 & & \\ & 4 & \\ 2 & & 1 \\ & & 3 \end{bmatrix}$$

In scipy:

matrix.indptr = [0, 1, 3, 5]

matrix.indices = [0, 1, 2, 0, 2]

matrix.data = [5, 4, 1, 2, 3]

Row 0 corresponds to the data from index 0 to 1 excluded

Compressed Sparse Row (CSR) format

$$\begin{bmatrix} 5 & & \\ & 4 & 1 \\ 2 & & 3 \end{bmatrix}$$

In scipy:

`matrix.indptr = [0, 1, 3, 5]`

`matrix.indices = [0, 1, 2, 0, 2]`

`matrix.data = [5, 4, 1, 2, 3]`

Row 1 corresponds to the data from index 1 to 3 excluded

© 2014 Pearson Education, Inc. or its affiliate(s). All rights reserved. Pearson Education, Inc., 501 Boylston Street, Boston, MA 02116

In scipy:

```
matrix.data = [ 5, 4, 1, 2, 3 ]
```

POLITECNICO MILANO 1863

Compressed Sparse Row (CSR) format

$$\begin{bmatrix} 5 & & \\ & 4 & 1 \\ 2 & & 3 \end{bmatrix}$$

In scipy:

```
matrix.indptr = [ 0, 1, 3, 5 ]  
matrix.indices = [ 0, 1, 2, 0, 2 ]  
matrix.data = [ 5, 4, 1, 2, 3 ]
```

Row 3 (index 5) does not exist in the matrix.

It is added because it allows to loop indptr without a special case for the last row

Compressed Sparse Row (CSR) format

Pros:

- Very fast row access (just two array accesses)
- Fast row-wise linear algebra operations

Cons:

- A bit less readable
- Very slow for column access, like COO

Compressed Sparse Column (CSC) format

The Compressed Sparse Column format works exactly like the CSR except that it compresses columns.

The *indptr* will represent the column starting point and the *indices* will represent the row.

How to use COO, CSC, CSR in practice

In practice, one usually creates the matrix using the simple COO format and then uses the library functions to get the CSC or CSR

Use the format you need at each step.

If first you need column access but later row access, it is likely better to change the format than to stick to a single one.

What happens if I use a CSR when I need fast column access, or a CSC for row access?

Nothing really...

the result will be correct (unless you use the internal data structures directly), but the script will be tens or hundreds of times slower