# DS-GA 1007
# Programming for Data Science

Lecture 12

pandas III + SQL II – Operations on Tables

Dask package in Python splits large tables into small tables and large tasks into dependent small tasks

# DS-GA 1007
# Programming for Data Science

Lecture 12

pandas III + SQL II – Operations on Tables

Dask package in Python allows us to scale pandas operations to large datasets

# DS-GA 1007
# Programming for Data Science

Lecture 12

pandas III + SQL II – Operations on Tables

Center for Data Science

NYU

# Announcements

- ▶ Homework 9 due **Monday November 25** at 11:59pm
- ▶ Project
  - ▶ Milestone due **Thursday November 28** at 11:59pm
  - ▶ Background and Plans
- ▶ Labs
  - ▶ Submit on Jupyter Hub under Assignments tab
  - ▶ Access scores from Submitted Assignments under Assignments tab

# Agenda

- ▶ Lesson
  - ▶ pandas
  - ▶ SQL
  - ▶ Dask
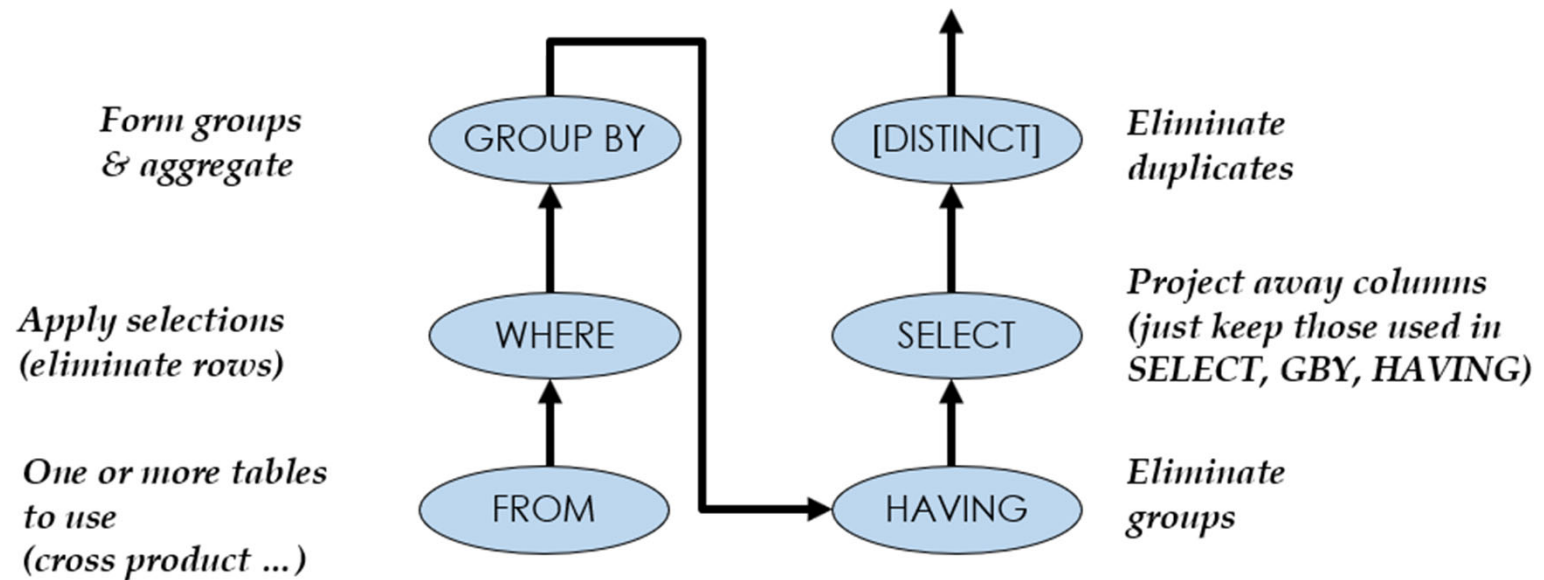- ▶ Demos
- ➡ Joins
  - ▶ Grouping
  - ▶ Pivot / Unpivot
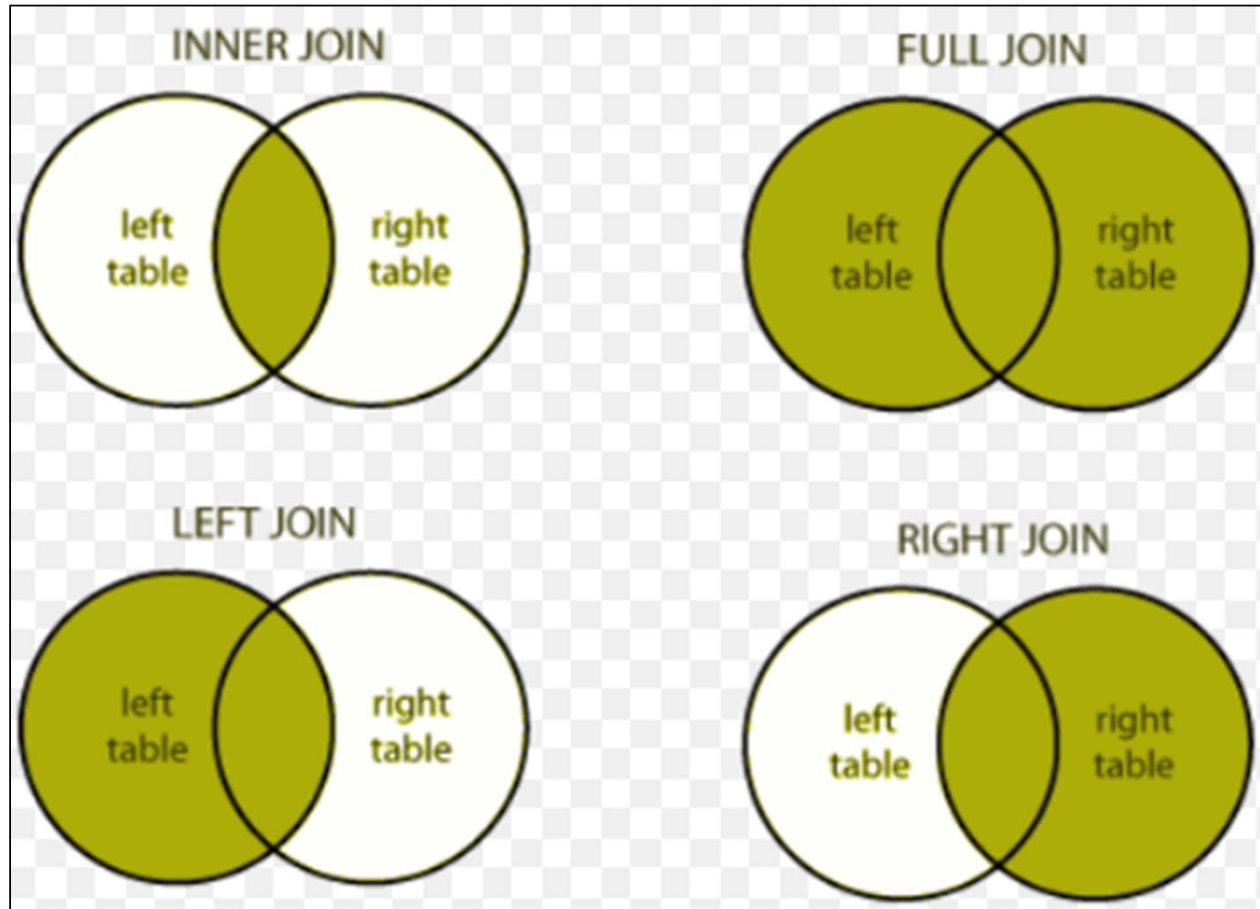  - ▶ Comparing running time and memory

## Objectives

- ▶ How can we combine multiple tables through joins?
- ▶ Group together rows by values in columns
- ▶ Pivot between rows and columns in a table
- ▶ How can we access large datasets?

# SQL Commands

```
SELECT      [DISTINCT] target-list
FROM        relation-list
WHERE       qualification
GROUP BY    grouping-list
HAVING      group-qualification
```



**Form groups & aggregate** → GROUP BY

**Apply selections (eliminate rows)** → WHERE

**One or more tables to use (cross product ...)** → FROM

**Eliminate duplicates** → [DISTINCT]

**Project away columns (just keep those used in SELECT, GBY, HAVING)** → SELECT

**Eliminate groups** → HAVING

# SQL Joins

# SQL Inner Join



```
SELECT * FROM s JOIN t ON s.u = t.v;

SELECT * FROM s INNER JOIN t ON s.u = t.v;

SELECT * FROM s, t WHERE s.u = t.v;
```

# SQL Left Outer Join



```
SELECT * FROM s LEFT JOIN t ON s.u = t.v;
```

# SQL Cross Join

## s

| M | U |
|---|---|
| 1 | W |
| 2 | X |
| 3 | X |
| 4 | Y |

## t

| N | V |
|---|---|
| A | X |
| B | X |
| C | Y |
| D | Z |

SELECT * FROM s, t;

| M | U | N | V |
|---|---|---|---|
| 1 | W | A | X |
| 2 | X | A | X |
| 3 | X | A | X |
| 4 | Y | A | X |
| 1 | W | B | X |
| 2 | X | B | X |
| 3 | X | B | X |
| 4 | Y | B | X |

(to be continued …)

(… continued)

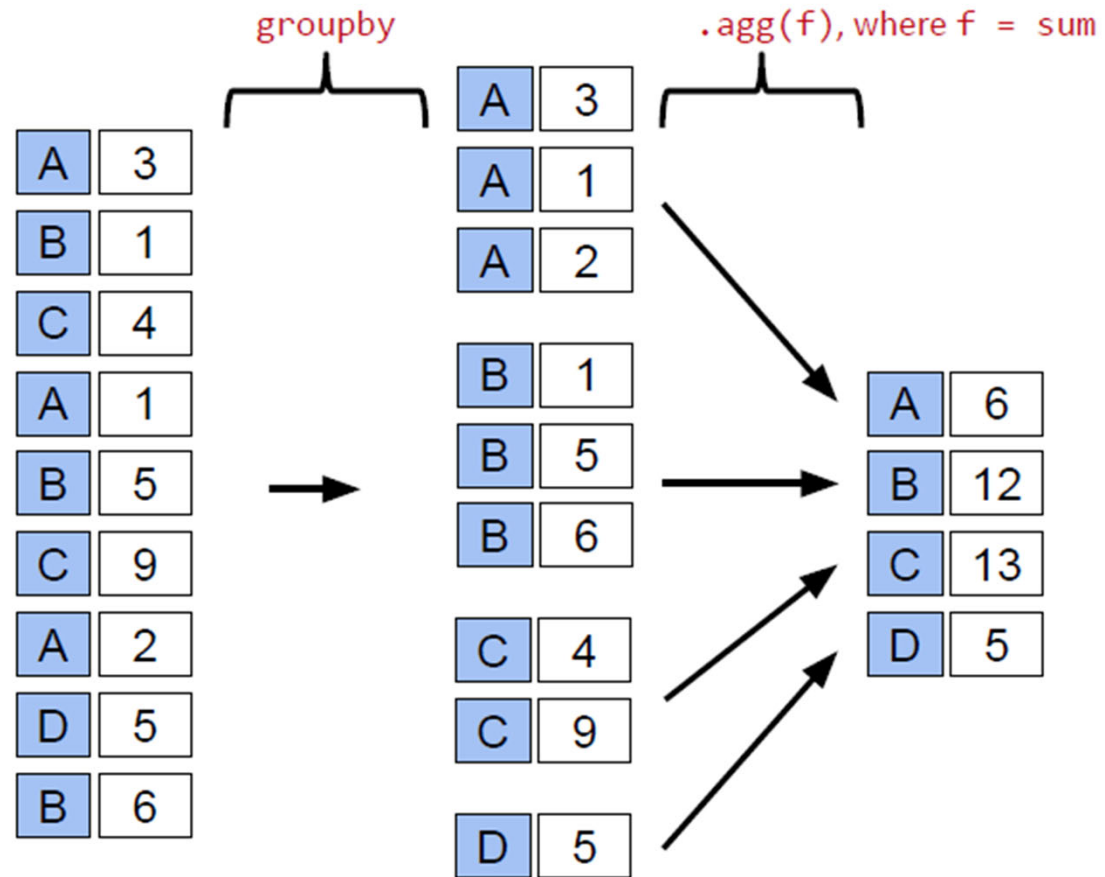| M | U | N | V |
|---|---|---|---|
| 1 | W | C | Y |
| 2 | X | C | Y |
| 3 | X | C | Y |
| 4 | Y | C | Y |
| 1 | W | D | Z |
| 2 | X | D | Z |
| 3 | X | D | Z |
| 4 | Y | D | Z |

# Agenda

- ▶ Lesson
  - ▶ pandas
  - ▶ SQL
  - ▶ Dask
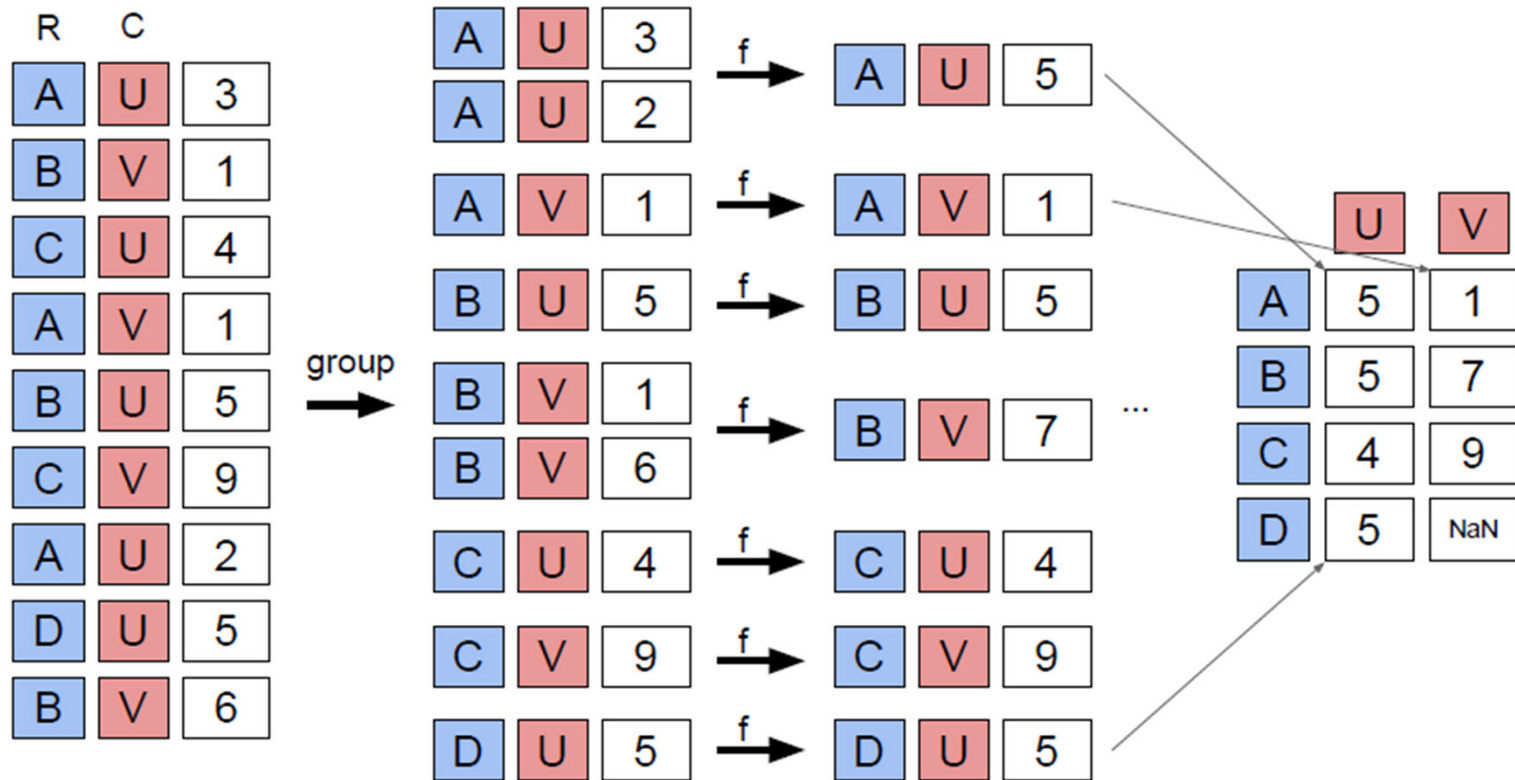- ▶ Demos
  - ▶ Joins
  - **→** Grouping
  - ▶ Pivot / Unpivot
  - ▶ Comparing running time and memory

## Objectives

- ▶ How can we combine multiple tables through joins?
- ▶ Group together rows by values in columns
- ▶ Pivot between rows and columns in a table
- ▶ How can we access large datasets?

# Group



groupby

.agg(f), where f = sum

# Agenda

- ▶ Lesson
  - ▶ pandas
  - ▶ SQL
  - ▶ Dask
- ▶ Demos
  - ▶ Joins
  - ▶ Grouping
  - → Pivot / Unpivot
  - ▶ Comparing running time and memory

## Objectives

- ▶ How can we combine multiple tables through joins?
- ▶ Group together rows by values in columns
- ▶ Pivot between rows and columns in a table
- ▶ How can we access large datasets?

# Pivot

# Agenda

▶ Lesson
  ▶ pandas
  ▶ SQL
  ▶ Dask
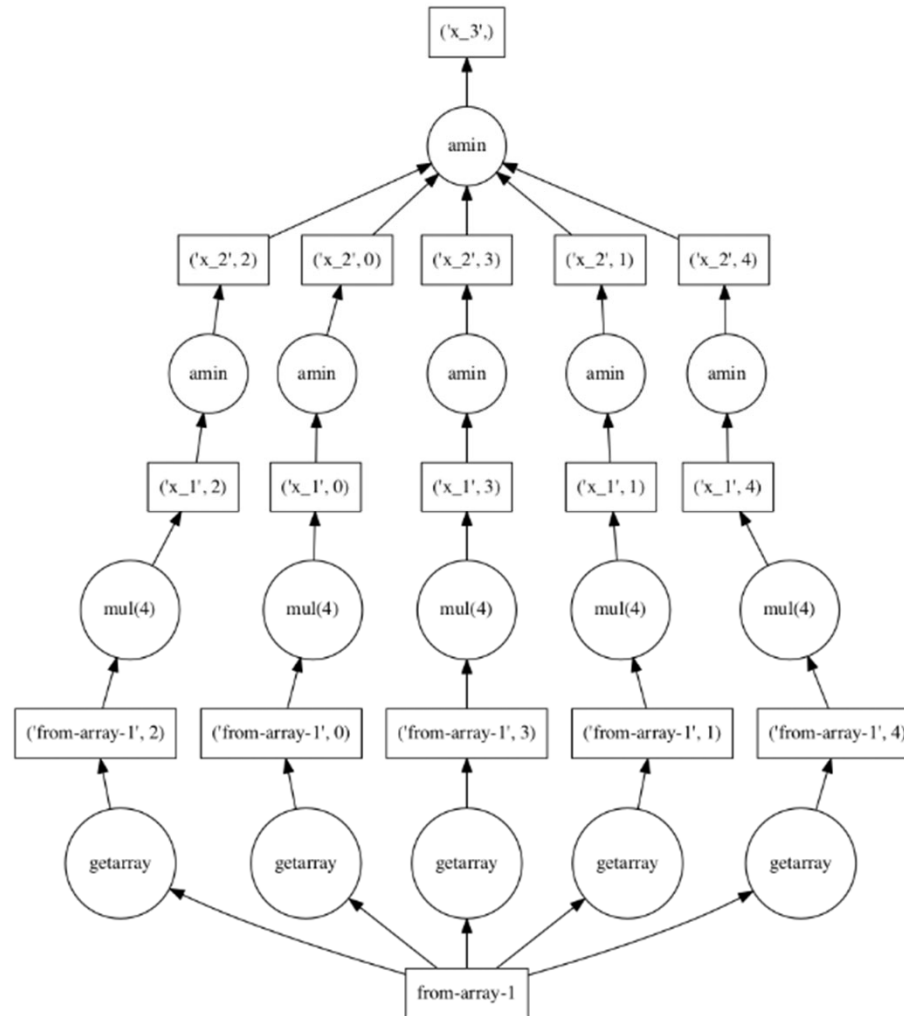▶ Demos
  ▶ Joins
  ▶ Grouping
  ▶ Pivot / Unpivot
  ▶ Comparing running time and memory

## Objectives

▶ How can we combine multiple tables through joins?

▶ Group together rows by values in columns

▶ Pivot between rows and columns in a table

▶ How can we access large datasets?

# Dask

# Summary

- pandas and SQL
  - Joining
  - Grouping
  - Pivot
- Dask
  - Running Time
  - Memory