Berkeley Engineering | BerkeleyHaas

**PROFESSIONAL CERTIFICATE IN MACHINE LEARNING AND ARTIFICIAL INTELLIGENCE**

**Module 4**
**Fundamentals of Data Analysis**
Office Hours with Viviana Márquez
September 21, 2023

- Slack
- Required activities for Module 4
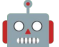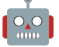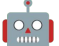- Content review Module 4: Fundamentals of Data Analytics
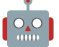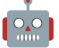- Questions

## Slack

#cohort-august-2023



*Slack Workspace
Invitation*

# Required Activities for Module 4

- 🤖 Codio Activity 4.1: BASIC JOINS ON DATASETS
- 🤖 Codio Activity 4.2: COMPLEX JOINS ON DATASETS
- 🤖 Codio Activity 4.3: CREATING SCATTERPLOTS, HISTOGRAMS, AND DISTRIBUTION PLOTS
- 🤖 Codio Activity 4.4: CREATING VIOLIN, BOX, AND JOINT PLOTS
- Try-It Activity 4.1: MORE SOPHISTICATED PLOTTING
- 🤖 Codio Activity 4.5: STRING OPERATIONS
- 🤖 Codoi Activity 4.6: DATA CLEANING
- Try-It Activity 4.2: ANALYZING A REAL-WORLD DATASET
- 🤖 Quiz 4.1: THE FUNDAMENTALS OF DATA ANALYSIS

## Content review Module 4: Fundamentals Data Analysis

- Joining tables
- Data cleaning
- Code
    - Pandas merge
    - Data cleaning
        - String operations

- Joining tables
- Data cleaning
- Code
  - Pandas merge
  - Data cleaning
    - String operations

# Joining tables

- In the real-world, often data will be split across multiple tables
- You will have to combine records from those tables based on related columns

# Types of joins

- **INNER JOIN**: Returns records that have matching values in both tables.
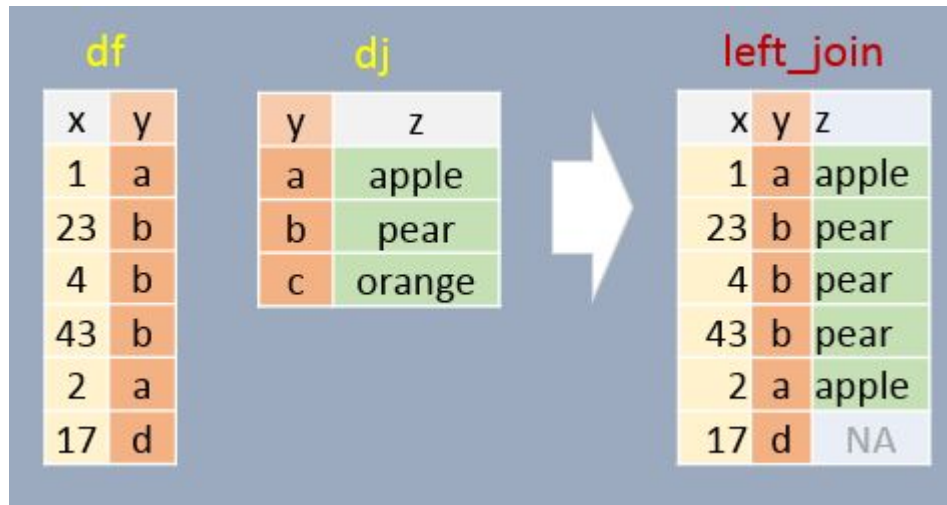- **LEFT (OUTER) JOIN**: Returns all records from the left table, and the matched records from the right table.
- **RIGHT (OUTER) JOIN**: Returns all records from the right table, and the matched records from the left table.
- **FULL (OUTER) JOIN**: Returns all records when there is a match in either the left or the right table.
- **CROSS JOIN**: Returns the Cartesian product of the two tables.
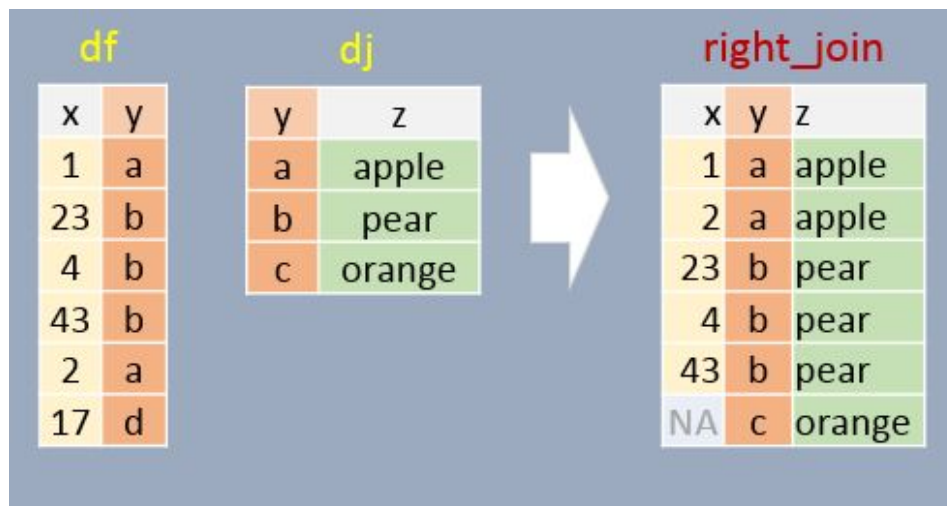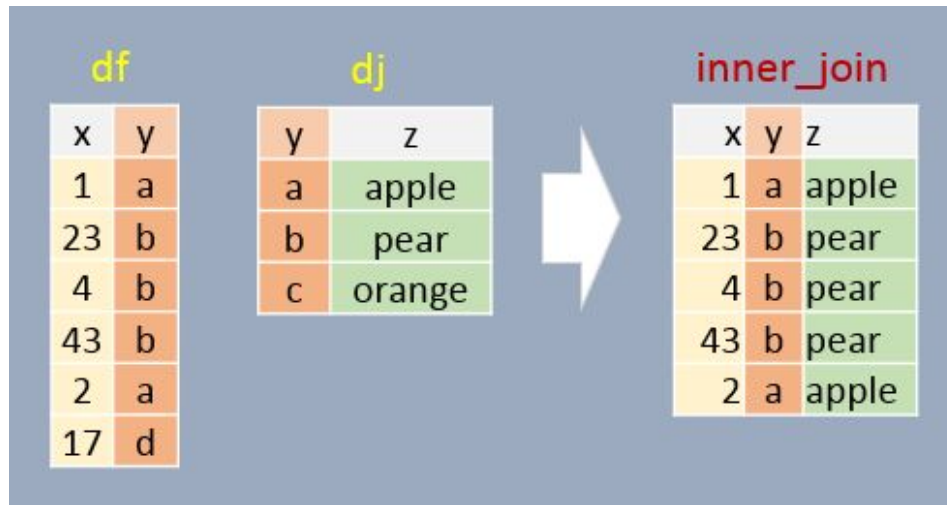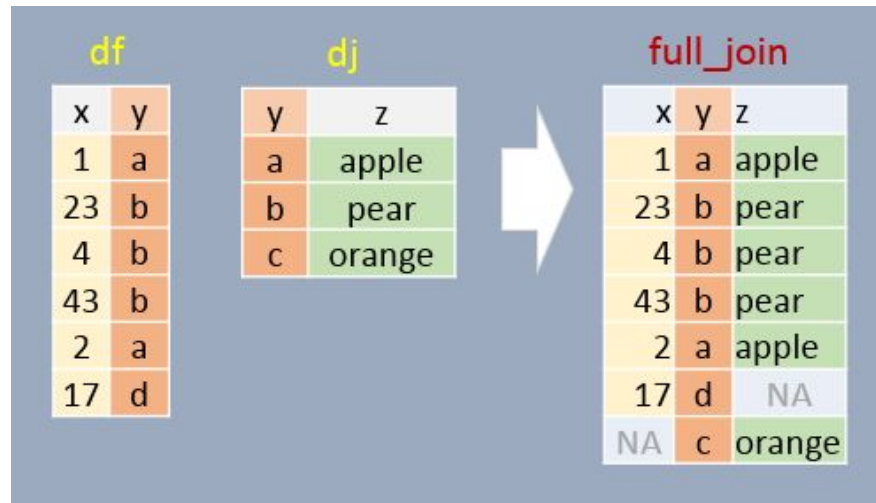- **SELF JOIN**: Joining a table with itself.

# Left join

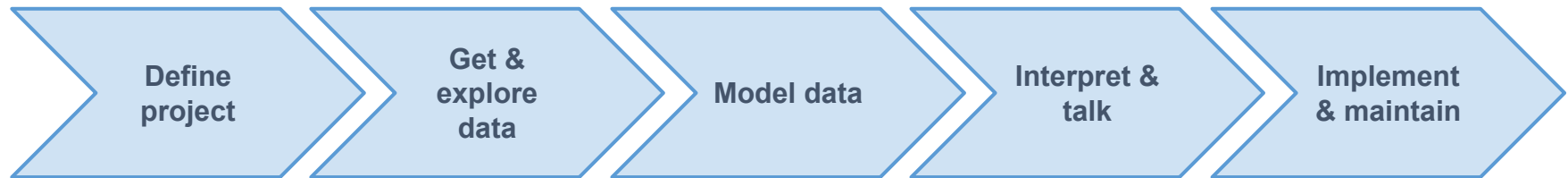# Right join

# Inner join

# Full join

# Let's code joins!

## Content review Module 4: Fundamentals Data Analysis

- ✅ Joining tables
- Data cleaning
- Code
  - Pandas merge
  - Data cleaning
    - String operations

# The Data Science Lifecycle

| Define project | Get & explore data | Model data | Interpret & talk | Implement & maintain |

## Define project

- Specify business problem
- Acquire domain knowledge

## Get and explore data

- Find appropriate data
- Exploratory Data Analysis
- Clean and pre-process data
- Feature engineering

## Model data

- Determine ML task
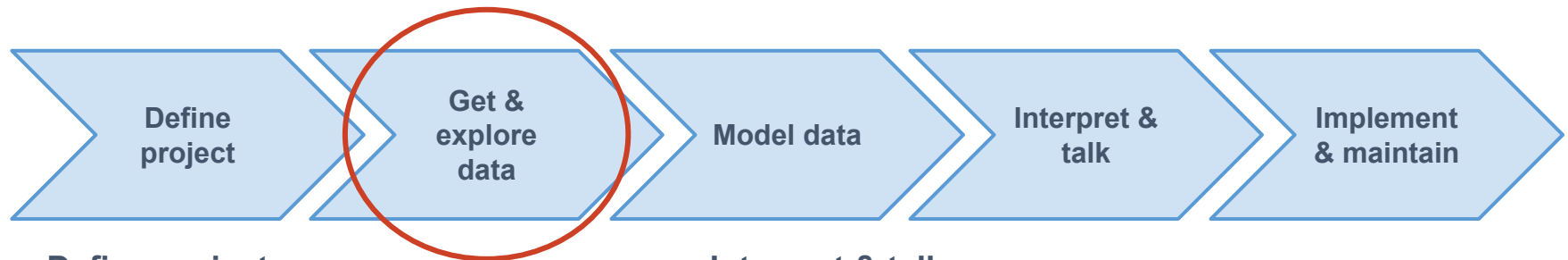- Build candidate models
- Select model based on performance metrics

## Interpret & talk

- Interpret model
- Communicate model insights

## Implement & maintain

- Set up function to predict on new data
- Document process
- Monitor and maintain model

# The Data Science Lifecycle

| Define project | Get & explore data | Model data | Interpret & talk | Implement & maintain |

**Define project**

- Specify business problem
- Acquire domain knowledge

**Get and explore data**

- Find appropriate data
- Exploratory Data Analysis
- Clean and pre-process data
- Feature engineering

**Model data**

- Determine ML task
- Build candidate models
- Select model based on performance metrics

**Interpret & talk**

- Interpret model
- Communicate model insights

**Implement & maintain**

- Set up function to predict on new data
- Document process
- Monitor and maintain model

# Data cleaning

- Identifying and correcting or removing errors, inaccuracies, and incomplete or irrelevant data
- Objective: Improve data quality, making it suitable for modeling

# Data cleaning

- Included but not limited to:
    - Handling missing values: drop them or account for them
    - Handling outliers: drop them or account for them or keep them
    - Remove duplicates
    - Handling incorrect data types
    - Handling inconsistent data (example: age shouldn't be negative)

- ✅ Joining tables
- ✅ Data cleaning
- Code
    - Pandas merge
    - Data cleaning
        - String operations

## Content review Module 4: Fundamentals Data Analysis

- ✅ Joining tables
- ✅ Data cleaning
- ✅ Code
    - Pandas merge
    - Data cleaning
        - String operations

# Let's code data cleaning!
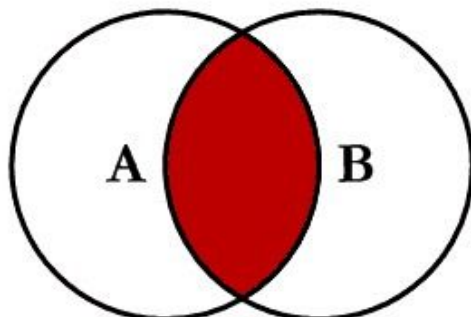
# QUESTIONS?

# Types of joins

# SQL JOINS
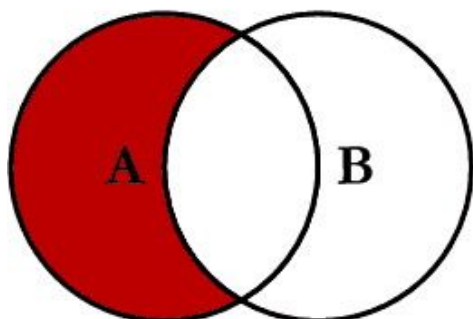


SELECT <select_list>
FROM TableA A
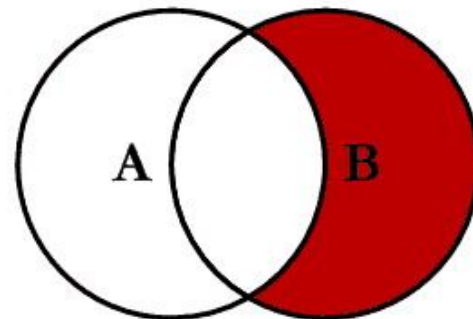LEFT JOIN TableB B
ON A.Key = B.Key

SELECT <select_list>
FROM TableA A
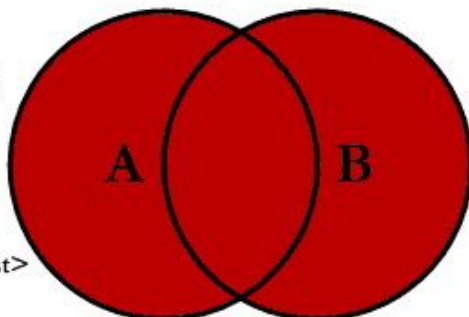RIGHT JOIN TableB B
ON A.Key = B.Key

SELECT <select_list>
FROM TableA A
INNER JOIN TableB B
ON A.Key = B.Key
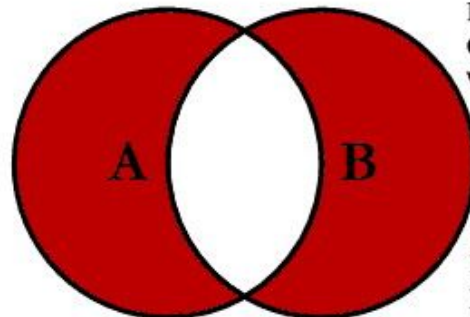
SELECT <select_list>
FROM TableA A
LEFT JOIN TableB B
ON A.Key = B.Key
WHERE B.Key IS NULL

SELECT <select_list>
FROM TableA A
RIGHT JOIN TableB B
ON A.Key = B.Key
WHERE A.Key IS NULL

SELECT <select_list>
FROM TableA A
FULL OUTER JOIN TableB B
ON A.Key = B.Key

SELECT <select_list>
FROM TableA A
FULL OUTER JOIN TableB B
ON A.Key = B.Key
WHERE A.Key IS NULL
OR B.Key IS NULL

© C.L. Moffatt, 2008