# YData: Introduction to Data Science



# Lecture 11: Joins

# Overview

Grouping continued

Pivot Tables

Joining tables

# Grouping

# Grouping by one column

The tb.group() method aggregates all rows with the same value for a column into a single row in the resulting table.
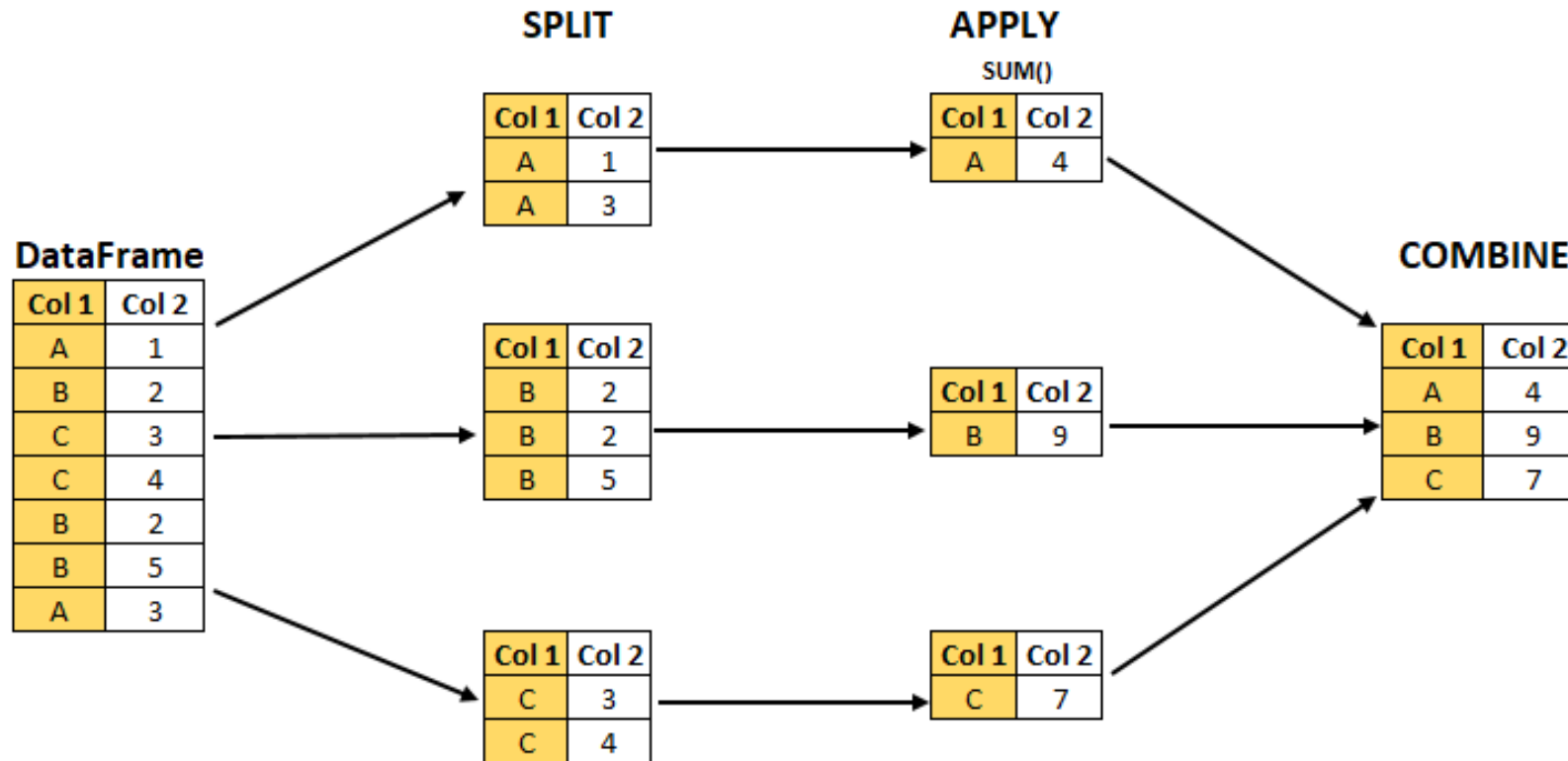
tb.group("grouping col", agg_function)

- "grouping col":  column to group data by
- agg_function:   function on how data in each group should be combined

Examples of aggregating functions:
- len: number of items in each group (default if no second argument is specified)
- list:  list of all values in each group
- sum: total of all grouped values

# Grouping: split-apply-combine



tb.group("Col 1", sum)

# Grouping by multiple columns

The tb.group() method can also aggregate values in rows that share the combination of values in multiple columns

tb.group(["grouping col1", "grouping col2"], agg_function)
- ["grouping col1", "grouping col2"]:  list of columns to group by
- agg_function:    function on how data in each group should be combined

Let's explore this in Jupyter!

# Pivot Tables

# Pivot Tables

Pivot tables aggregate values according to two categorical variables but the results are in a <u>table</u>

- i.e., same as grouping by two categorical variables but puts one variable as the rows and the other as columns

Produces a grid of counts or aggregated values two required arguments:

- First: variable that forms column labels of grid
- Second: variable that forms row labels of grid

Two optional arguments (include both or neither)

- values = `column label to aggregate'
- collect = function with which to aggregate

**Let's explore this in Jupyter!**

**Grouping**

cat var 1          cat var 2

| Flavor | Color | count |
|---|---|---|
| bubblegum | pink | 1 |
| chocolate | dark brown | 2 |
| chocolate | light brown | 1 |
| strawberry | pink | 2 |

**Pivot Table**

cat var 1

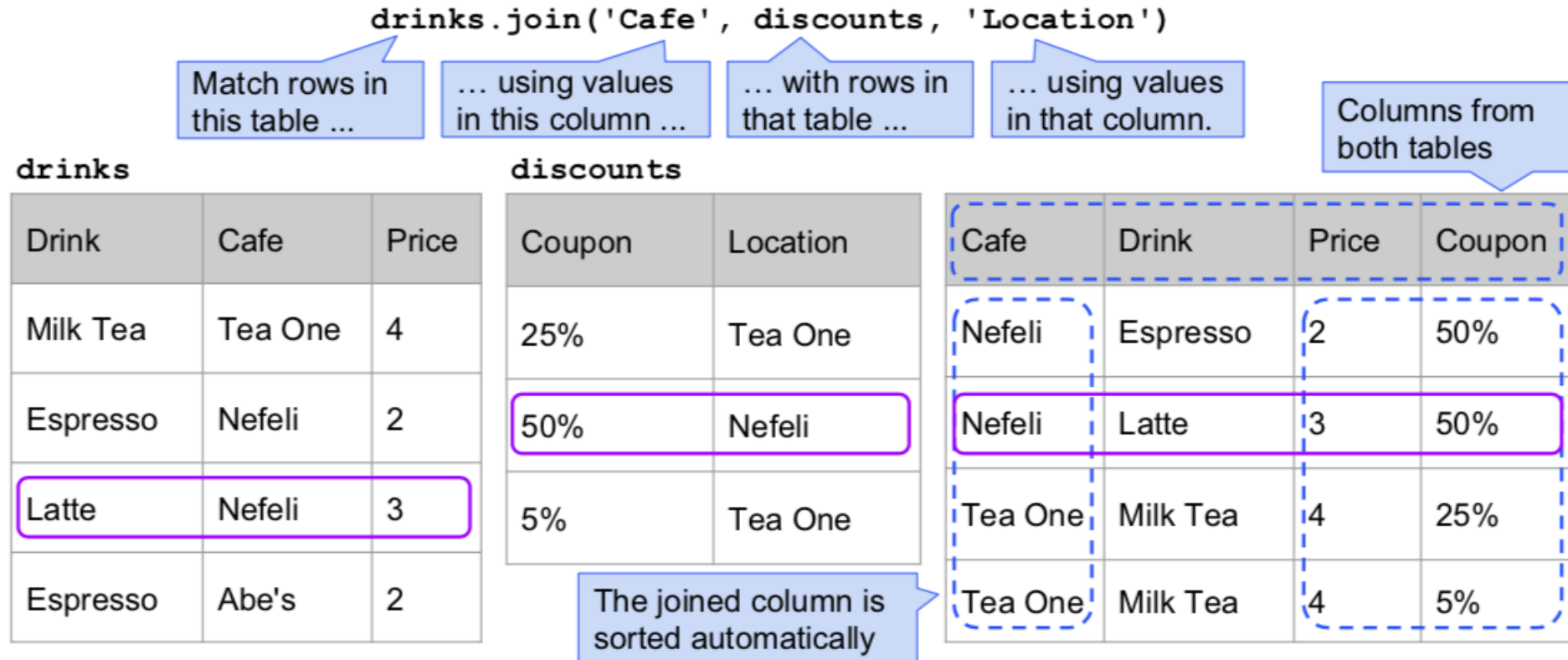|  | Color | bubblegum | chocolate | strawberry |
|---|---|---|---|---|
| cat var 2 | dark brown | 0 | 2 | 0 |
|  | light brown | 0 | 1 | 0 |
|  | pink | 1 | 0 | 2 |

# Joins

# Joining Two Tables

Joining involves combining the rows of two tables together into a new table
- A column in each table needs to be specified which indicates how the rows should be combined

tb1.join("col tb1", tb2, "col tb2")
- tb1: the first table
- "col tb1":   a column in the first table
- tb2:   the second table
- "col tb2":   a column in the second table

# Joining Two Tables



Let's explore this in Jupyter!