

Grouping

Learn how DataFrames can be grouped based on particular columns.

Chapter Goals:

- Learn how to group DataFrames by columns
- Write code to retrieve home run statistics through DataFrame grouping

A. Grouping by column

When dealing with large amounts of data, it is usually a good idea to group the data by common categories. For example, we could group a large dataset of MLB player statistics by year, so we can deal with each year's data separately.

With pandas DataFrames, we can perform dataset grouping with the `groupby` function. A common usage of the function is to group a DataFrame by values from a particular column, e.g. a column representing years.

The code below shows how to use the `groupby` function, with the example of grouping MLB data by year.

```
# Predefined df of MLB stats
print('{}\n'.format(df))

groups = df.groupby('yearID')
for name, group in groups:
    print('Year: {}'.format(name))
    print('{}\n'.format(group))

print('{}\n'.format(groups.get_group(2016)))
print('{}\n'.format(groups.sum()))
print('{}\n'.format(groups.mean()))
```



The grouping code example produced three DataFrames for the years 2015, 2016, and 2017. The three DataFrame groups are contained in the `groups`

variable, and we used its `sum` and `mean` functions to retrieve the total and average per-year statistics.

In addition to aggregation functions like `sum` and `mean`, we can also filter the groups using `filter`. The `filter` function takes in another function as its required argument, which specifies how we want to filter the groups. The output of `filter` is the concatenation of all the groups that pass the filter.

The code below shows how to use the `filter` function.

```
no2015 = groups.filter(lambda x: x.name > 2015)
print(no2015)
```



In the above code example, the lambda function passed into `filter` returns `True` if the group (represented as `x`) represents a year greater than 2015. The output is the concatenation of the 2016 and 2017 groups.

B. Multiple columns

DataFrame grouping is not just limited to a single column. Rather than passing a single column label into `groupby`, we can use a list of column labels to specify grouping by multiple columns.

Grouping by multiple columns can be useful when multiple data features have many different values. For example, if our dataset consisted of MLB players, grouping by both team and year would give us an organized way to view a team's roster throughout the years.

```
# player_df is predefined
groups = player_df.groupby(['yearID', 'teamID'])

for name, group in groups:
    print('Year, Team: {}'.format(name))
    print('{}\n'.format(group))

print(groups.sum())
```



In the code above, we grouped the MLB data by both year and team, resulting

in each group's name being a tuple of year and team. Using the `sum` function, we obtained the annual total hits for each team.

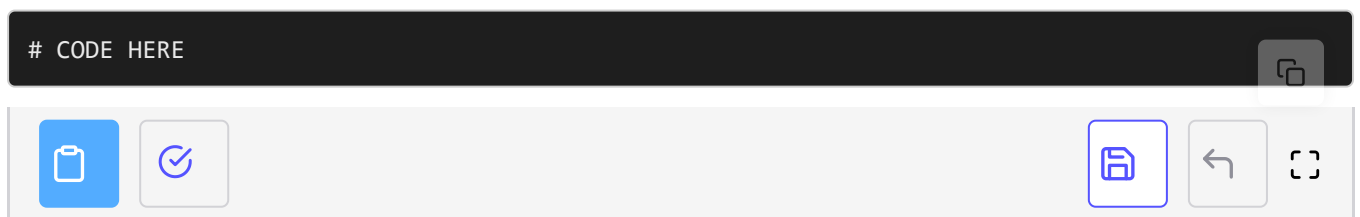
Time to Code!

The coding exercises for this chapter involve performing grouping operations on `df`, which contains all MLB batting data from 1871-2017. Using `df`, our goal is to retrieve home run (HR) statistics for 2017.

To do this, we need to calculate the total number of home runs hit each year. This involves first grouping `df` by year.

Set `year_group` equal to `df.groupby` with `'yearID'` as the lone argument.

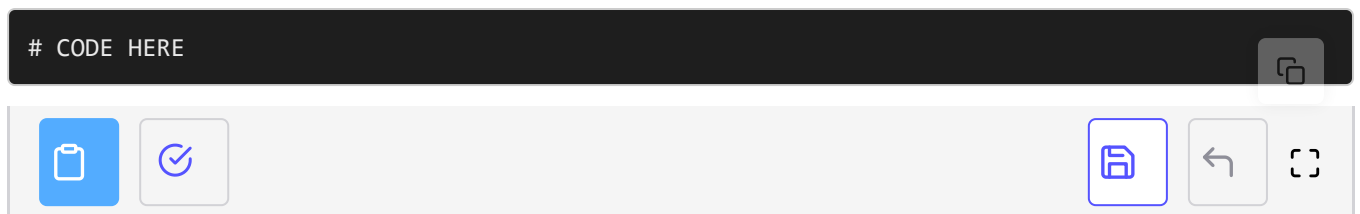
```
# CODE HERE
```



The yearly stats can be obtained from summing the values across the year-separated groups.

Set `year_stats` equal to `year_group.sum` applied with no arguments.

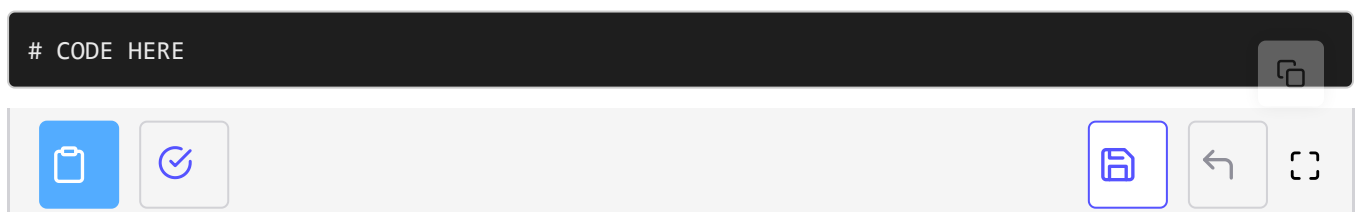
```
# CODE HERE
```



The `year_stats` DataFrame represents the total value for each stat per year. The row labels are the years and the column labels are the stat categories, e.g. home runs. Using the `loc` property, we'll retrieve the home run total for 2017.

Set `hr_2017` equal to `year_stats.loc` with `2017` as the row index and `'HR'` as the column index.

```
# CODE HERE
```

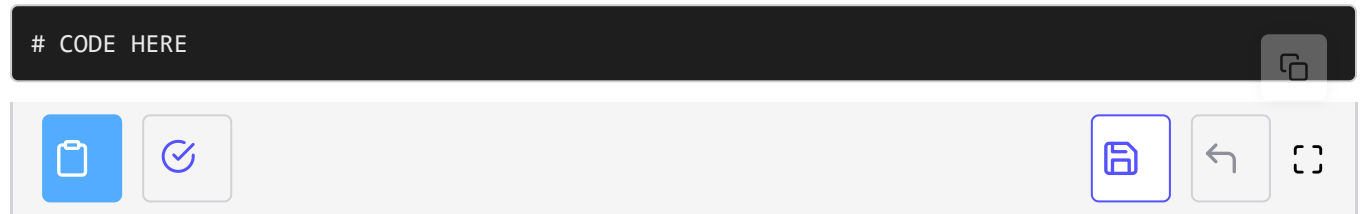


Next we want to get the yearly totals for each batting statistic *per team*. To do

this, we group the data by both the year and team.

Set `year_team_group` equal to `df.groupby` applied with the list `['yearID', 'teamID']`.

```
# CODE HERE
```



Once again, to obtain the yearly stats we just sum over all the groups.

Set `year_team_stats` equal to `year_team_group.sum` applied with no arguments.

```
# CODE HERE
```

