

DataFrame

Learn about the pandas DataFrame object for 2-D data.

Chapter Goals:

- Learn about the pandas DataFrame object and its basic utilities
- Write code to create and manipulate a pandas DataFrame

A. 2-D data

One of the main purposes of pandas is to deal with tabular data, i.e. data that comes from tables or spreadsheets. Since tabular data contains rows and columns, it is 2-D. For working with 2-D data, we use the

`pandas.DataFrame` (<https://pandas.pydata.org/pandas-docs/stable/reference/api/pandas.DataFrame.html>) object, which we'll refer to simply as a DataFrame.

A DataFrame is created through the `pd.DataFrame` constructor, which takes in essentially the same arguments as `pd.Series`. However, while a Series could be constructed from a scalar (representing a single value Series), a DataFrame cannot.

Furthermore, `pd.DataFrame` takes in an additional `columns` keyword argument, which represents the labels for the columns (similar to how `index` represents the row labels).

The code below shows how to use the `pd.DataFrame` constructor.

```
1 df = pd.DataFrame()
2 # Newline added to separate DataFrames
3 print('{}\n'.format(df))
4
5 df = pd.DataFrame([5, 6])
6 print('{}\n'.format(df))
7
8 df = pd.DataFrame([[5,6]])
9 print('{}\n'.format(df))
10
11 df = pd.DataFrame([[5, 6], [1, 3]],
```



```

12             index=['r1', 'r2'],
13             columns=['c1', 'c2'])
14 print('{}\n'.format(df))
15
16 df = pd.DataFrame({'c1': [1, 2], 'c2': [3, 4]},
17                   index=['r1', 'r2'])
18 print('{}\n'.format(df))

```



Output

0.934s

```

Empty DataFrame
Columns: []
Index: []

```

```

      0
0  5
1  6

      0  1

```

Note that when we use a Python dictionary for initialization, the DataFrame takes the dictionary's keys as its column labels.

B. Upcasting

When we initialize a DataFrame of mixed types, upcasting occurs on a per-column basis. The `dtypes` property returns the types in each column as a Series of types.

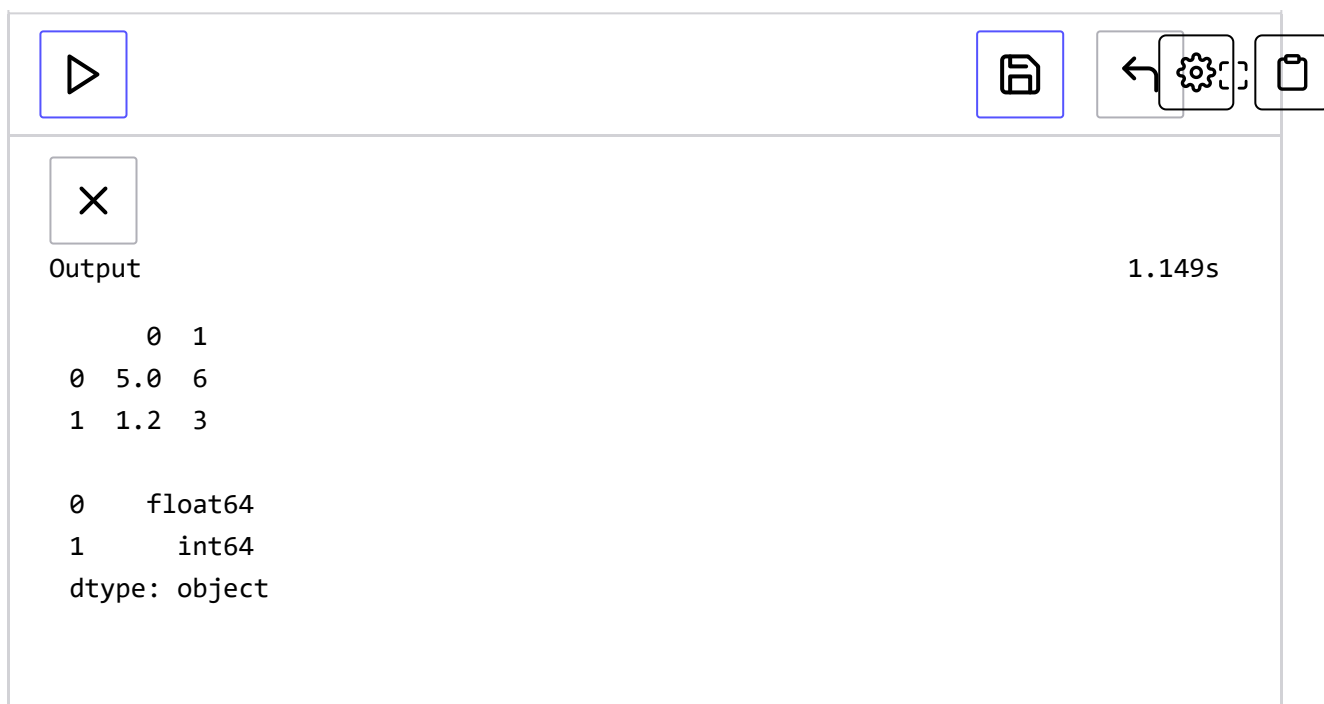
The code below shows how upcasting works in DataFrames. You'll notice that upcasting only occurs in the first column for the DataFrame below, because the second column's values are both integers.

```

1 upcast = pd.DataFrame([[5, 6], [1.2, 3]])
2 print('{}\n'.format(upcast))
3 # Datatypes of each column
4 print(upcast.dtypes)

```





Output 1.149s

```
   0  1
0  5.0  6
1  1.2  3

0    float64
1     int64
dtype: object
```

C. Appending rows

We can append additional rows to a given DataFrame through the `append` (<https://pandas.pydata.org/pandas-docs/stable/reference/api/pandas.DataFrame.append.html>) function. The required argument for the function is either a Series or DataFrame, representing the row(s) we append.

Note that the `append` function returns the modified DataFrame but doesn't actually change the original. Furthermore, when we append a Series to the DataFrame, we either need to specify the `name` for the series or use the `ignore_index` keyword argument. Setting `ignore_index=True` will change the row labels to integer indexes.

The code below shows example usages of the `append` function.

```
1 df = pd.DataFrame([[5, 6], [1.2, 3]])
2 ser = pd.Series([0, 0], name='r3')
3
4 df_app = df.append(ser)
5 print('{}\n'.format(df_app))
6
7 df_app = df.append(ser, ignore_index=True)
8 print('{}\n'.format(df_app))
9
10 df2 = pd.DataFrame([[0,0],[9,9]])
11 df_app = df.append(df2)
12 print('{}\n'.format(df_app))
13
```

The screenshot shows a Jupyter Notebook interface. At the top right, there are icons for settings (gear) and a clipboard. Below these, on the left, is a play button icon, and on the right, are icons for saving (floppy disk), undo (curved arrow), and a zoom icon. The main area contains a code cell with a close button (X) in the top left. The code cell's output is displayed below it, showing two DataFrames. The first DataFrame has columns 0, 1, and 2, with rows indexed 0, 1, and r3. The second DataFrame has the same columns but includes an additional row indexed 2. A vertical scrollbar is on the right side of the output area, and a timing indicator '0.770s' is shown in the top right corner of the output area.

```

0  1
0  5.0  6
1  1.2  3
r3  0.0  0

0  1
0  5.0  6
1  1.2  3
2  0.0  0

```

D. Dropping data

We can drop rows or columns from a given DataFrame through the `drop` (<https://pandas.pydata.org/pandas-docs/stable/reference/api/pandas.DataFrame.drop.html>) function. There is no required argument, but the keyword arguments of the function gives us two ways to drop rows/columns from a DataFrame.

The first way is using the `labels` keyword argument to specify the labels of the rows/columns we want to drop. We use this alongside the `axis` keyword argument (which has default value of `0`) to drop from the rows or columns axis.

The second method is to directly use the `index` or `columns` keyword arguments to specify the labels of the rows or columns directly, without needing to use `axis`.

The code below shows examples on how to use the `drop` function.

```

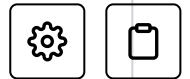
1 df = pd.DataFrame({'c1': [1, 2], 'c2': [3, 4],
2                   'c3': [5, 6]},
3                   index=['r1', 'r2'])
4 # Drop row r1
5 df_drop = df.drop(labels='r1')

```

```

5 df_drop = df.drop(labels='c1',
6 print('{}\n'.format(df_drop))
7
8 # Drop columns c1, c3
9 df_drop = df.drop(labels=['c1', 'c3'], axis=1)
10 print('{}\n'.format(df_drop))
11
12 df_drop = df.drop(index='r2')
13 print('{}\n'.format(df_drop))
14
15 df_drop = df.drop(columns='c2')
16 print('{}\n'.format(df_drop))
17
18 df.drop(index='r2', columns='c2')
19 print('{}\n'.format(df_drop))

```



Output

0.783s

```

      c1  c3
r1     1   5
r2     2   6

```

```

      c1  c3
r1     1   5
r2     2   6

```



Similar to `append`, the `drop` function returns the modified DataFrame but doesn't actually change the original.

Note that when using `labels` and `axis`, we can't drop both rows and columns from the DataFrame.

Time to Code!

The coding exercise for this chapter involves creating various pandas DataFrame objects.

We'll first create a DataFrame from a Python dictionary. The dictionary will have key-value pairs 'c1':[0, 1, 2, 3] and 'c2':[5, 6, 7, 8], in that order.

The index for the DataFrame will come from the list of row labels ['r1', 'r2', 'r3', 'r4'] .

Set df equal to pd.DataFrame with the specified dictionary as the first argument and the list of row labels as the index keyword argument.

```
1 df = pd.DataFrame({'c1': [0, 1, 2, 3], 'c2': [5, 6, 7, 8]},
2                     index=['r1', 'r2', 'r3', 'r4'])
3
4
5
```

Show Results

Show Console

×

1 of 1 Tests Passed

Result	Input	Expected Output	Actual Output	Reason
✓		c1 c2 r1 0 5 r2 1 6 r3 2 ...	c1 c2 r1 0 5 r2 1 6 r3 2 ...	Value of df is correct, good job!

1.397s

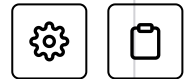
We'll create another DataFrame, this one representing a single row. Rather than a dictionary for the first argument, we use a list of lists, and manually set the column labels to ['c1', 'c2'] .

Since there is only one row, the row labels will be ['r5'] .

Set row_df equal to pd.DataFrame with [[9, 9]] as the first argument, and the specified column and row labels for the columns and index keyword arguments.

```
1 row_df = pd.DataFrame([[9, 9]], columns=['c1', 'c2'], index=['r5'])
```

2
3
4



Show Results

Show Console



1 of 1 Tests Passed

Result	Input	Expected Output	Actual Output	Reason
✓		c1 c2 r5 9 9	c1 c2 r5 9 9	Value of row_df is correct, good job!

0.608s

After creating `row_df`, we append it to the end of `df` and drop row `'r2'`.

Set `df_app` equal to `df.append` with `row_df` as the only argument.

Then set `df_drop` equal to `df_app.drop` with `'r2'` as the `labels` keyword argument.

```
1 df_app = df.append(row_df)
2 df_drop = df_app.drop(labels='r2')
3
```



Show Results

Show Console



2 of 2 Tests Passed

Result	Input	Expected Output	Actual Output	Reason
✓		c1 c2 r1 0 5 r2 1 6 r3 2 "	c1 c2 r1 0 5 r2 1 6 r3 2 "	Value of <code>df_app</code> is correct, good job!

Result	Input	Expected Output	Actual Output	Reason
✓		c1 c2 r1 0 5 r3 2 7 r4 4 ...	c1 c2 r1 0 5 r3 2 7 r4 4 ...	Value of df_drop is correct, good job!

0.885s

← Back

Series

Next →

Combining

✓ Mark as Completed

16% completed, meet the criteria and claim your course certificate!

Buy Certificate

ⓘ Report an Issue

🔗 Ask a Question
(https://discuss.educative.io/tag/dataframe__data-analysis-with-pandas__machine-learning-for-software-engineers)