

**Python** Data Analysis



DANIEL Y. CHEN

# Pandas for Everyone

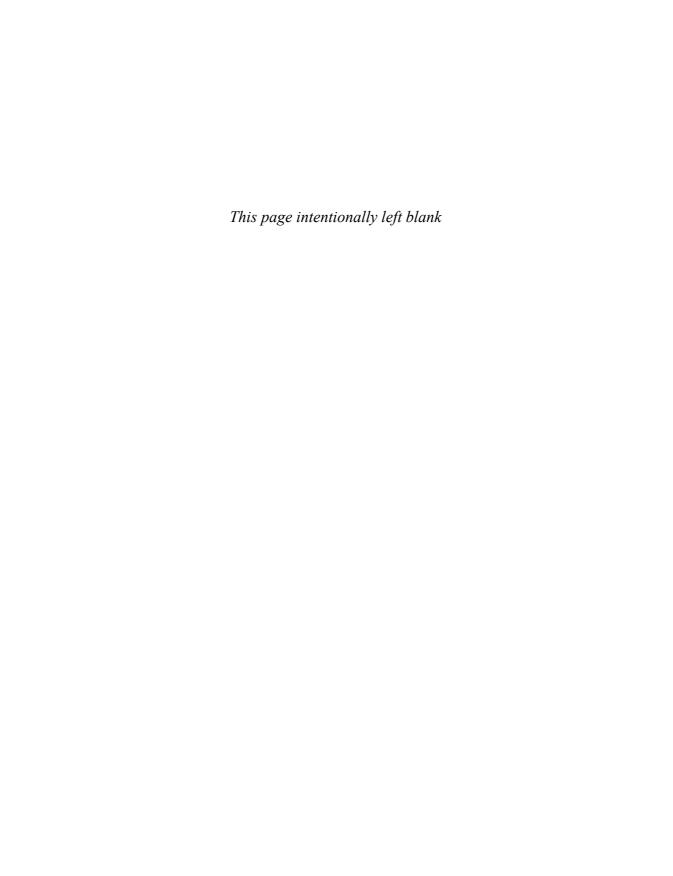
# Part II

# Data Manipulation

**Chapter 4** Data Assembly

Chapter 5 Missing Data

Chapter 6 Tidy Data



## Data Assembly

#### 4.1 Introduction

By now, you should be able to load data into Pandas and do some basic visualizations. This part of the book focuses on various data cleaning tasks. We begin with assembling a data set for analysis by combining various data sets together.

## **Concept Map**

- 1. Prior knowledge
  - a. loading data
  - b. subsetting data
  - c. functions and class methods

## **Objectives**

This chapter will cover:

- 1. Tidy data
- 2. Concatenating data
- 3. Merging data sets

## 4.2 Tidy Data

Hadley Wickham, <sup>1</sup> one of the more prominent members of the R community, talks about the idea of *tidy* data. In fact, he's written a paper about this concept in the *Journal of Statistical Software*. <sup>2</sup> Tidy data is a framework to structure data sets so they can be easily analyzed. It is mainly used as a goal one should aim for when cleaning data. Once you understand what tidy data is, that knowledge will make data collection much easier.

<sup>1.</sup> Hadley Wickham's homepage: http://hadley.nz

<sup>2.</sup> Tidy data paper: http://vita.had.co.nz/papers/tidy-data.pdf

So what is *tidy* data? Hadley Wickham's paper defines it as meeting the following criteria:

- Each row is an observation.
- Each column is a variable.
- Each type of observational unit forms a table.

#### 4.2.1 Combining Data Sets

We begin with Hadley Wickham's last tidy data point: "Each type of observational unit forms a table." When data is tidy, you need to combine various tables together to answer a question. For example, there may be a separate table holding company information and another table holding stock prices. If we want to look at all the stock prices within the tech industry, we may first have to find all the tech companies from the company information table, and then combine that data with the stock price data to get the data we need for our question. The data may have been split up into separate tables to reduce the amount of redundant information (we don't need to store the company information with each stock price entry), but this arrangement means we as data analysts must combine the relevant data ourselves to answer our question.

At other times, a single data set may be split into multiple parts. For example, with time-series data, each date may be in a separate file. In another case, a file may have been split into parts to make the individual files smaller. You may also need to combine data from multiple sources to answer a question (e.g., combine latitudes and longitudes with zip codes). In both cases, you will need to combine data into a single dataframe for analysis.

#### 4.3 Concatenation

One of the (conceptually) easier ways to combine data is with concatenation. Concatenation can be thought of appending a row or column to your data. This approach is possible if your data was split into parts or if you performed a calculation that you want to append to your existing data set.

Concatenation is accomplished by using the concat function from Pandas.

#### 4.3.1 Adding Rows

Let's begin with some example data sets so you can see what is actually happening.

```
import pandas as pd

df1 = pd.read_csv('../data/concat_1.csv')
df2 = pd.read_csv('../data/concat_2.csv')
df3 = pd.read_csv('../data/concat_3.csv')

print(df1)

A B C D
0 a0 b0 c0 d0
1 a1 b1 c1 d1
2 a2 b2 c2 d2
3 a3 b3 c3 d3
```

```
print(df2)
             С
                 D
        b4
            c4 d4
   а5
        b5
           c5 d5
        b6
            c6
                d6
3
                d7
   а7
        b7
            с7
print(df3)
                С
                      D
          b8
                     d8
     a8
               с8
1
     а9
          b9
                     d9
               с9
         b10
                    d10
   a10
              c10
         b11
                    d11
   a11
              c11
```

Stacking the dataframes on top of each other uses the concat function in Pandas. All of the dataframes to be concatenated are passed in a list.

```
row_concat = pd.concat([df1, df2, df3])
print(row_concat)
      Α
           В
                 С
                       D
                     d0
0
     a0
          b0
                c0
                     d1
1
     a1
          b1
                с1
2
     a2
          b2
                c2
                     d2
3
     a3
          b3
                c3
                     d3
0
     a4
          b4
                c4
                     d4
 1
     а5
          b5
                c5
                     d5
2
     a6
          b6
                c6
                     d6
3
     а7
          b7
                с7
                     d7
0
     а8
          b8
                с8
                     d8
     a9
          b9
                c9
                     d9
2
   a10
         b10
               c10
                    d10
         b11
               c11
                    d11
   a11
```

As you can see, concat blindly stacks the dataframes together. If you look at the row names (i.e., the row indices), they are also simply a stacked version of the original row indices. If we apply the various subsetting methods from Table 2.3, the table will be subsetted as expected.

#### Question

What happens when you use loc or ix to subset the new dataframe?

Section 2.2.1 showed the process for creating a Series. However, if we create a new series to append to a dataframe, it does not append correctly.

```
# create a new row of data
new_row_series = pd.Series(['n1', 'n2', 'n3', 'n4'])
print(new_row_series)
     n1
     n2
     n3
3
     n4
dtype: object
# attempt to add the new row to a dataframe
print(pd.concat([df1, new_row_series]))
               С
          В
                    D
              c0
0
    a0
         b0
                   d0
                       NaN
1
         b1
                   d1 NaN
    a1
              c1
2
    a2
         b2 c2
                   d2 NaN
3
    а3
         b3 c3 d3 NaN
   NaN NaN NaN NaN
                        n1
   NaN NaN NaN NaN
                        n2
2
        NaN
             NaN
                  NaN
                        n3
   NaN
        NaN
             NaN
                  NaN
                        n4
```

The first things you may notice are the NaN values. This is simply Python's way of representing a "missing value" (see Chapter 5, "Missing Data"). We were hoping to append our new values as a row, but that didn't happen. In fact, not only did our code not append the values as a row, but it also created a new column completely misaligned with everything else.

If we pause to think about what is happening here, we can see that the results actually make sense. First, if we look at the new indices that were added, we notice that they are very similar to the results we obtained when we concatenated dataframes earlier. The indices of the new\_row series object are analogs to the row numbers of the dataframe. Also, since our series did not have a matching column, our new\_row was added to a new column.

To fix this problem, we can turn our series into a dataframe. This data frame contains one row of data, and the column names are the ones the data will bind to.