

▼ Data Wrangling

David Gerbing
The School of Business
Portland State University
gerbing@pdx.edu

Data in the real world typically does not arrive ready for analysis. Instead usually manipulate the data in various ways to derive a nice, clean, tidy data frame of rows by columns with all the data values in a column of the same type, such as character strings, integers, or floating point numbers (i.e., with decimal digits).

■ **Data Wrangling:** The process of cleaning, tidying, and otherwise preparing data for analysis.

The following examples demonstrate some useful data manipulations that are applicable to all data analysis: subsetting a data table, converting variable types, and variable transformations. Merging data frames, another common data manipulation, is shown elsewhere.

This task of data wrangling is where most data scientists spend most of their time, up to 80% is a common understanding. Larger organizations have added a new job category called data engineer, a specialist in data wrangling and other aspects of handling data, such as data collection.

The material in this and related notebooks covers some of the basic, and most common, data wrangling procedures.

▼ Preliminaries

▼ Packages

```
from datetime import datetime as dt
now = dt.now()
print ("Analysis on", now.strftime("%Y-%m-%d"), "at", now.strftime("%H:%M"))

Analysis on 2021-06-27 at 14:06
```

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
```

▼ Read

```
#d = pd.read_excel('data/employee.xlsx')
d = pd.read_excel('http://lessRstats.com/data/employee.xlsx')
```

```
d.shape
```

```
(37, 9)
```

```
d.head()
```

	Name	Years	Gender	Dept	Salary	JobSat	Plan	Pre	Post
0	Ritchie, Darnell	7.0	M	ADMN	53788.26	med	1	82	92
1	Wu, James	NaN	M	SALE	94494.58	low	1	62	74
2	Hoang, Binh	15.0	M	SALE	111074.86	low	3	96	97

▼ Data Frame Row Identifiers

Columns in a data frame represent the variables in the analysis. Columns can be identified by the corresponding variable name or by their numerical position in the data frame.

■ **Index:** An integer that specifies the numerical row or column position.

Unfortunately (in my opinion), Python begins all counting with 0 instead of 1. However, the principle of identifying a row or a column by the corresponding integer remains.

In the `d` data frame that contains the read data, as shown above, the default row identifiers are the row indices. However, if there is column of unique identifiers in the data table, that column can be designated as the row identifiers.

In the above `d` data table, the `Name` column appears as any other variable in the data frame. However, `Name` is not a variable per se to analyze but an ID field, with a unique value for each row. Replace the default integer row labels with the values of the column `Name` with the `set_index()` function.

Note that data manipulation methods typically do not change the original data frame. To save changes, explicitly save the manipulation into a variable, such as below, where the change to the `d` data frame is saved back into the `d` data frame. If there is no assigned variable for the output, the output is directed to the console. You will be able to view the output, but no changes are saved.

```
d = d.set_index('Name')
d.head()

      Years Gender Dept  Salary JobSat Plan Pre Post
Name
Ritchie, Darnell  7.0   M  ADMN  53788.26  med   1   82   92
Wu, James  NaN   M  SALE  94494.58  low   1   62   74
Hoang, Binh  15.0   M  SALE  111074.86  low   3   96   97
Jones, Alissa  5.0   F  NaN  53772.58  NaN   1   65   62
```

Or, set the row names as the data values directly when read with a function such as `read_excel()`. Use the parameter `index_col=0` to set the column index. The variable `Name` in the original data frame is in the first column, that is, Column 0.

```
d = pd.read_excel('http://lessRstats.com/data/employee.xlsx', index_col=0)
#d = pd.read_excel('data/employee.xlsx', index_col=0)
d.head()
```

	Years	Gender	Dept	Salary	JobSat	Plan	Pre	Post
Name								
Ritchie, Darnell	7.0	M	ADMN	53788.26	med	1	82	92
Wu, James	NaN	M	SALE	94494.58	low	1	62	74
Hoang, Binh	15.0	M	SALE	111074.86	low	3	96	97
Jones, Alissa	5.0	F	NaN	53772.58	NaN	1	65	62

▼ Subset Rows and Columns

A data frame is the Python representation of a rectangular data table with rows and columns. In many situations, we wish to view or analyze just a portion of the entire data frame, a subset, an extraction from the original. Many possibilities of subsetting exist: a single data value (cell), a single row, a single column, and a range of rows or columns.

Express all references to data values in a data frame in terms of those rows and columns. Reference data values within data frame `d` by its rows and columns:

```
d[row_reference, column_reference]
```

The reference can be a corresponding name of a row or column or its associated integer index. If referencing location by name, use the `.loc()` function for "location". Reference the integer index with the `.iloc()` function for "integer location" or "index location". The "weird" part, as with all Python counting, is that counting rows or columns starts with 0 instead of starting at 1.

Note: There is no assignment of the sub-setted information to the new data frame in the examples below, just a display of the requested information.

▼ Select a Single Cell

We see from the original data table that Binh Hoang's salary is \$111,074.86.

To display the data value stored in a single, specific cell, specify a single row reference and a single column reference. Here reference the location of the data value by the row name and column name, so use the `loc` method for "location".

Note that the output is not assigned to another data frame, so the output is directed to the space right below where the cell is located.

```
d.loc['Hoang, Binh', 'Salary']
```

```
111074.86
```

Or, reference the location of the data value by its row and column index according to their integer positions in the data frame with `iloc`. The data for Binh Hoang are in the third row, Row 2 because Python starts counting at 0. The variable `Salary` is in the fourth column, Column 3.

```
d.iloc[2,3]
```

```
111074.86
```

▼ Select Multiple Cells

A colon, `:`, indicates a range of either rows (before the comma) or columns (after the comma).

```
d2 = d.loc['Wu, James':'Jones, Alissa', 'Gender':'JobSat']
d2.head()
```

	Gender	Dept	Salary	JobSat	Plan	Pre	Post
Name							
Wu, James	M	SALE	94494.58	low	1	62	74
Hoang, Binh	M	SALE	111074.86	low	3	96	97
Jones, Alissa	F	NaN	53772.58	NaN	1	65	62

When subsetting row or column indices in pandas with the colon operator, `:`, always reference the counting system that begins with 0. And then reference the row after the last row you wish to select.

For example, to select the second through the fourth row, the pandas row indices are 1 through 3. So specify a row range of 1:4. To select the second through fourth row, specify a column range of 1:5. (Why do the Python people have to complicate such a simple issue as counting?)

To illustrate, show here again the first five rows of the `d` data frame.

```
d.head()
```

	Years	Gender	Dept	Salary	JobSat	Plan	Pre	Post
Name								
Ritchie, Darnell	7.0	M	ADMN	53788.26	med	1	82	92
Wu, James	NaN	M	SALE	94494.58	low	1	62	74
Hoang, Binh	15.0	M	SALE	111074.86	low	3	96	97
Jones, Alissa	5.0	F	NaN	53772.58	NaN	1	65	62

Now select the second through the fourth rows, getting three rows indicated by indices 1 through 3, and the second through the fifth columns, getting four columns, indicated by indices 1 through 4.

```
d2 = d.iloc[1:4, 1:5]
```

```
d2
```

	Gender	Dept	Salary	JobSat	Plan	Pre	Post
Name							
Wu, James	M	SALE	94494.58	low	1	62	74
Hoang, Binh	M	SALE	111074.86	low	3	96	97
Jones, Alissa	F	NaN	53772.58	NaN	1	65	62

▼ Subset Rows

▼ Select a Single Row

The optional, but good practice to include.

```
d.loc['Hoang, Binh', 'Salary']
```

```
111074.86
```

▼ Select Multiple Rows

A colon, `:`, indicates a range of either rows (before the comma) or columns (after the comma).

```
d2 = d.loc['Wu, James':'Jones, Alissa', 'Gender':'JobSat']
```

```
d2.head()
```

	Gender	Dept	Salary	JobSat	Plan	Pre	Post
Name							
Wu, James	M	SALE	94494.58	low	1	62	74
Hoang, Binh	M	SALE	111074.86	low	3	96	97
Jones, Alissa	F	NaN	53772.58	NaN	1	65	62

Here, `query()` function may be more straightforward, but there is another expression for the extraction that many people use. This expression involves repeating the name of the data frame, which becomes awkward for long data frame names and a logical expression that involves multiple variables. Still, this form is frequently encountered in pandas data manipulation.

```
dd = d[d['Salary'] > 100000]
```

```
dd
```

	Years	Gender	Dept	Salary	JobSat	Plan	Pre	Post
Name								
Hoang, Binh	15.0	M	SALE	111074.86	low	3	96	97

▼ Subset Rows by Data Values

Retrieve just those rows of data that match a logical condition. Here identify all rows of data for employees with a salary more than \$100,000 per year. The first example only displays the result. The second example creates a new data frame with the result, and then displays the subset.

This first example invokes the most straightforward method of selecting rows that satisfy a logical condition: The `query()` function. Note the logical condition is enclosed in quotes.

```
d.query('Salary > 100000')
```

```
dd
```

	Years	Gender	Dept	Salary	JobSat	Plan	Pre	Post
Name								
Hoang, Binh	15.0	M	SALE	111074.86	low	3	96	97

The `query()` function may be more straightforward, but there is another expression for the extraction that many people use. This expression involves repeating the name of the data frame, which becomes awkward for long data frame names and a logical expression that involves multiple variables. Still, this form is frequently encountered in pandas data manipulation.

```
dd = d[d['Salary'] > 100000]
```

```
dd
```

	Years	Gender	Dept	Salary	JobSat	Plan	Pre	Post
Name								
Hoang, Binh	15.0	M	SALE	111074.86	low	3	96	97

Here, `query()` function may be more straightforward, but there is another expression for the extraction that many people use. This expression involves repeating the name of the data frame, which becomes awkward for long data frame names and a logical expression that involves multiple variables. Still, this form is frequently encountered in pandas data manipulation.

```
dd = d[d['Salary'] > 100000]
```

```
dd
```

	Years	Gender	Dept	Salary	JobSat	Plan	Pre	Post
Name								
Hoang, Binh	15.0	M	SALE	111074.86	low	3	96	97

▼ Subset Rows by Data Values

Retrieve just those rows of data that match a logical condition. Here identify all rows of data for employees with a salary more than \$100,000 per year. The first example only displays the result. The second example creates a new data frame with the result, and

```
d2 = d.iloc[:, [1,3]]
```

```
d2.head()
```

Gender Salary

Name	Gender	Salary
Ritchie, Darnell	M	53788.26
Wu, James	M	94494.58
Hoang, Binh	M	111074.86
Jones, Alissa	F	53772.58
Downs, Deborah	F	57139.90

▼ Select by Variable Type

The select function for selecting by variable types is `select_dtypes()`, either with the parameter `exclude` or `include`.

```
d_num = d.select_dtypes(exclude=['object'])
```

```
d_num.head()
```

Years Salary Plan Pre Post

Name	Years	Salary	Plan	Pre	Post
Ritchie, Darnell	7.0	53788.26	1	82	92
Wu, James	NaN	94494.58	1	62	74
Hoang, Binh	15.0	111074.86	3	96	97
Jones, Alissa	5.0	53772.58	1	65	62
Downs, Deborah	7.0	57139.90	2	90	86

```
d_obj = d.select_dtypes(include=['object'])
```

```
d_obj.head()
```

Gender Dept JobSat

Name	Gender	Dept	JobSat
Ritchie, Darnell	M	ADMN	med
Wu, James	M	SALE	low
Hoang, Binh	M	SALE	low
Jones, Alissa	F	NaN	NaN
Downs, Deborah	F	FINC	high

▼ Chained Functions

Likely the most straightforward to accomplish data frame subsets is with the `query()` and `filter()` functions. However, when doing multiple function calls, one after the other, you can chain these calls into a single call. This chaining elucidates a complex, multi-step data manipulation process with highly readable, structured code.

To specify a set of chained functions, include the entire expression within parentheses `()`, then separate each function call on its own line. In this example, subset by rows, then by columns, then sort on the `Salary` column with the `sort_values()` function.

```
(d
  .query('Salary > 10000')
  .filter(['Gender', 'Salary'])
  .sort_values(['Salary'], ascending=False)
)
```

Gender Salary

Name	Gender	Salary
Correll, Trevor	M	134419.23
James, Leslie	F	122563.38
Hoang, Binh	M	111074.86
Capelle, Adam	M	108138.43

▼ Delete Rows or Columns

Delete a Row

Start with 37 rows of data.

```
d.shape
```

```
(37, 8)
```

Use the `drop()` function to delete a row by row name, to result in 36 rows.

```
d2 = d.drop(['Wu, James'])
```

```
d2.shape
```

```
(36, 8)
```

Delete a row by row index, here the second row, to result in 36 rows.

```
d2 = d.drop([d.index[1]])
```

```
d2.shape
```

```
(36, 8)
```

The data for James Wu is gone.

```
d2.head()
```

Years Gender Dept Salary JobSat Plan Pre Post

Name	Years	Gender	Dept	Salary	JobSat	Plan	Pre	Post
Ritchie, Darnell	7.0	M	ADMN	53788.26	med	1	82	92
Hoang, Binh	15.0	M	SALE	111074.86	low	3	96	97
Jones, Alissa	5.0	F	NaN	53772.58	NaN	1	65	62
Downs, Deborah	7.0	F	FINC	57139.90	high	2	90	86
Afshari, Anbar	6.0	F	ADMN	69441.93	high	2	100	100

▼ Delete a Column

The function `drop()` deletes a row or column from the data frame, as specified by the `axis` parameter. The default value of `axis` is `'rows'`, so if dropping a column, need to explicitly specify. Here drop the variable `Plan` from the resulting data frame of only numeric variables because it is an integer coded categorical variable. Dropping that variable leaves only continuous variables.

```
d_num = d2.drop(['Plan'], axis='columns')
```

```
d_num.head()
```

Years Gender Dept Salary JobSat Plan Pre Post

Name	Years	Gender	Dept	Salary	JobSat	Plan	Pre	Post
Ritchie, Darnell	7.0	M	ADMN	53788.26	med	1	82	92
Hoang, Binh	15.0	M	SALE	111074.86	low	3	96	97
Jones, Alissa	5.0	F	NaN	53772.58	NaN	1	65	62
Downs, Deborah	7.0	F	FINC	57139.90	high	2	90	86
Afshari, Anbar	6.0	F	ADMN	69441.93	high	2	100	100

▼ Other Issues

A kind of "weird" issue exists (especially if you are used to R). Subsets (slices) of data frames by default do not create clean copies. The new data frame is still linked to the original data frame. Changes to the newly created data frame change the original as well! Modifications to the data or indices of the copy are reflected back to the original object.

To make the subset data frame completely independent of `d`, what pandas calls a `deep copy`, invoke the subset extraction with the `copy()` function. This function is not needed if there will not be further manipulation of the contents of the newly created data frame. Unless memory is an issue for very large data sets, or you know that there will be no further modification, using `copy()` is a good practice.

```
d2 = d.loc[:, 'Salary'].copy()
```

```
d2.head()
```

```
Name: Salary, dtype: float64
```

A short-hand specification to select a column does not call any method or function and only includes the names of the relevant columns. When learning a language, better to focus on the more complete implementation of a concept. However, there is a need to know this abbreviated form because it appears often in real-world applications.

```
d2 = d['Salary']
```

```
d2.head()
```

```
Name: Salary, dtype: float64
```

▼ Variable Transformation

Transform the values of a numerical variable with an equation, a function that specifies how each value is to be transformed. Transform the values of a categorical variable with a recoding that specifies how each individual value is to be replaced with a new value.

Numeric Variable

To transform the values of a numeric variable is straightforward: Enter the corresponding equation that defines the transformation. Refer to a variable within a data frame with the data frame name, such as `d`, and then the variable name within quotes and square brackets. For example, the refer to the `Salary` variable in the `d` data frame: `d['Salary']`.

This example creates a new variable, `Salary000`, defined as the original `Salary` variable with values divided by 1000, rounded to two decimal digits. The new variable is added to the already existing variables in the `d` data frame. Running the equation creates the values of the new variable for all rows of data in the data frame.

```
d['Salary000'] = round(d['Salary'] / 1000, 2)
```

```
d.head()
```

Years Gender Dept Salary JobSat Plan Pre Post Salary000

Name	Years	Gender	Dept	Salary	JobSat	Plan	Pre	Post	Salary000
Ritchie, Darnell	7.0	M	ADMN	53788.26	med	1	82	92	53.79
Hoang, Binh	15.0	M	SALE	111074.86	low	3	96	97	111.07
Jones, Alissa	5.0	F	NaN	53772.58	NaN	1	65	62	53.77
Downs, Deborah	7.0	F	FINC	57139.90	high	2	90	86	57.14
Afshari, Anbar	6.0	F	ADMN	69441.93	high	2	100	100	69.44

Categorical Variable

The function `replace()` replaces individual values of a categorical variable. Parameter `to_replace` indicates the values to be replaced, and `value` indicates the replacement value. Here replace across the entire data frame. Here replace both values of `Gender` with one statement.

```
d.dtypes
```

```
Years          float64
Gender         object
Salary         float64
JobSat         float64
Plan           int64
Salary000      float64
```

```
dtype: object
```

```
dtypes: float64
```

```
object
```

```
int64
```

```
float64
```

```
float64</pre
```