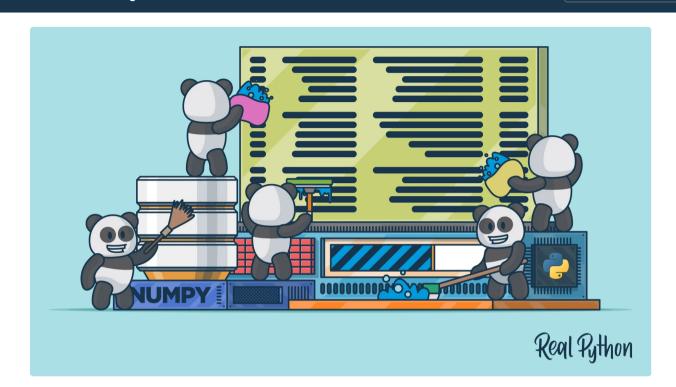
Real Rython

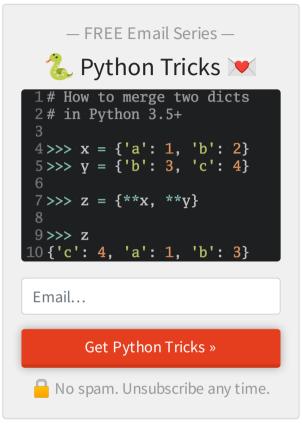


► Learn Python ▼

Pythonic Data Cleaning With pandas and NumPy

Changing the Index of a DataFrame

• Tidying up Fields in the Data



All Tutorial Topics advanced api basics best-practices community databases data-science devops django docker flask front-end gamedev gui intermediate machine-learning projects python testing tools web-dev web-scraping



- Combining str Methods with NumPy to Clean Columns
- Cleaning the Entire Dataset Using the applymap Function
- Renaming Columns and Skipping Rows
- Python Data Cleaning: Recap and Resources

THE WORLD'S MOST ELABORATE TRADING CHALLENGE





Remove ads

Watch Now This tutorial has a related video course created by the Real Python team. Watch it together with the written tutorial to deepen your understanding: **Data**Cleaning With pandas and NumPy

Data scientists spend a large amount of their time cleaning datasets and getting them down to a form with which they can work. In fact, a lot of data scientists argue that the initial steps of obtaining and cleaning data constitute 80% of the job.

Therefore, if you are just stepping into this field or planning to step into this field, it is important to be able to deal with messy data, whether that means missing values, inconsistent formatting, malformed records, or nonsensical outliers.

In this tutorial, we'll leverage Python's pandas and NumPy libraries to clean data.

We'll cover the following:

- Dropping unnecessary columns in a DataFrame
- Changing the index of a DataFrame
- Using .str() methods to clean columns
- Using the DataFrame.applymap() function to clean the entire dataset, element-wise
- Renaming columns to a more recognizable set of labels



Table of Contents

- Dropping Columns in a DataFrame
- Changing the Index of a DataFrame
- Tidying up Fields in the Data
- Combining str Methods with NumPy to Clean Columns
- Cleaning the Entire Dataset Using the applymap Function
- Renaming Columns and Skipping Rows
- Python Data Cleaning: Recap and Resources



• Skipping unnecessary rows in a CSV file

Free Bonus: Click here to get access to a free NumPy Resources Guide that points you to the best tutorials, videos, and books for improving your NumPy skills.

Here are the datasets that we will be using:

- BL-Flickr-Images-Book.csv A CSV file containing information about books from the British Library
- university_towns.txt A text file containing names of college towns in every US state
- olympics.csv A CSV file summarizing the participation of all countries in the Summer and Winter Olympics

You can download the datasets from Real Python's GitHub repository in order to follow the examples here.

Note: I recommend using Jupyter Notebooks to follow along.

This tutorial assumes a basic understanding of the pandas and NumPy libraries, including Panda's workhorse Series and DataFrame objects, common methods that can be applied to these objects, and familiarity with NumPy's NaN values.

Let's import the required modules and get started!





Dranning Calumne in a Data Frama

DIOPPING COLUMNIS III a Dataframe

Often, you'll find that not all the categories of data in a dataset are useful to you. For example, you might have a dataset containing student information (name, grade, standard, parents' names, and address) but want to focus on analyzing student grades.

In this case, the address or parents' names categories are not important to you. Retaining these unneeded categories will take up unnecessary space and potentially also bog down runtime.

pandas provides a handy way of removing unwanted columns or rows from a DataFrame with the drop() function. Let's look at a simple example where we drop a number of columns from a DataFrame.

First, let's create a DataFrame out of the CSV file 'BL-Flickr-Images-Book.csv'. In the examples below, we pass a relative path to pd.read_csv, meaning that all of the datasets are in a folder named Datasets in our current working directory:

```
Python
                                                                             >>>
>>> df = pd.read_csv('Datasets/BL-Flickr-Images-Book.csv')
>>> df.head()
    Identifier
                            Edition Statement
                                                    Place of Publication \
          206
0
                                         NaN
                                                                 London
1
          216
                                              London; Virtue & Yorston
                                         NaN
2
          218
                                         NaN
                                                                 London
          472
3
                                         NaN
                                                                 London
               A new edition, revised, etc.
          480
                                                                 London
  Date of Publication
                                    Publisher \
          1879 [1878]
                             S. Tinsley & Co.
0
                 1868
                                 Virtue & Co.
1
                 1869
                       Bradbury, Evans & Co.
                 1851
                                James Darling
                 1957 Warthaim & Macintoch
```

```
MEL LIIETIII & MACTILLOSII
                                                         Author \
                                               Title
0
                  Walter Forbes. [A novel.] By A. A
                                                          Α. Α.
1 All for Greed. [A novel. The dedication signed... A., A. A.
  Love the Avenger. By the author of "All for Gr... A., A. A.
  Welsh Sketches, chiefly ecclesiastical, to the... A., E. S.
   [The World in which I live, and my place in it... A., E. S.
                                   Contributors Corporate Author \
0
                                FORBES, Walter.
                                                              NaN
  BLAZE DE BURY, Marie Pauline Rose - Baroness
                                                              NaN
  BLAZE DE BURY, Marie Pauline Rose - Baroness
                                                              NaN
                   Appleyard, Ernest Silvanus.
3
                                                              NaN
                            BROOME, John Henry.
                                                              NaN
   Corporate Contributors Former owner
                                        Engraver Issuance type \
0
                                                   monographic
                      NaN
                                   NaN
                                             NaN
                                                   monographic
1
                      NaN
                                   NaN
                                             NaN
                                                   monographic
                      NaN
                                   NaN
                                             NaN
                                                   monographic
                      NaN
                                   NaN
                                             NaN
                                                   monographic
                      NaN
                                   NaN
                                             NaN
                                          Flickr URL \
0 http://www.flickr.com/photos/britishlibrary/ta...
 http://www.flickr.com/photos/britishlibrary/ta...
  http://www.flickr.com/photos/britishlibrary/ta...
  http://www.flickr.com/photos/britishlibrary/ta...
  http://www.flickr.com/photos/britishlibrary/ta...
                            Shelfmarks
    British Library HMNTS 12641.b.30.
0
    British Library HMNTS 12626.cc.2.
1
    British Library HMNTS 12625.dd.1.
    British Library HMNTS 10369.bbb.15.
    British Library HMNTS 9007.d.28.
```

When we look at the first five entries using the head() method, we can see that a handful of columns provide ancillary information that would be helpful to the library but isn't very

descriptive of the books themselves: Edition Statement, Corporate Author,
Corporate Contributors, Former owner, Engraver, Issuance type and Shelfmarks.

We can drop these columns in the following way:

Above, we defined a list that contains the names of all the columns we want to drop. Next, we call the drop() function on our object, passing in the inplace parameter as True and the axis parameter as 1. This tells pandas that we want the changes to be made directly in our object and that it should look for the values to be dropped in the columns of the object.

When we inspect the DataFrame again, we'll see that the unwanted columns have been removed:

```
Python
                                                                          >>>
>>> df.head()
  Identifier
                   Place of Publication Date of Publication \
0
          206
                                 London
                                               1879 [1878]
              London; Virtue & Yorston
          216
                                                       1868
          218
                                London
                                                      1869
         472
                                London
                                                       1851
                                London
         480
                                                       1857
               Publisher
                                                                      Title \
       S. Tinslev & Co.
                                         Walter Forbes. [A novel.] By A. A
0
           Virtue & Co. All for Greed. [A novel. The dedication signed...
  Bradbury, Evans & Co. Love the Avenger. By the author of "All for Gr...
          James Darling Welsh Sketches, chiefly ecclesiastical, to the...
3
   Wertheim & Macintosh [The World in which I live, and my place in it...
                                                     Flickr URL
     Author
      A. A. http://www.flickr.com/photos/britishlibrary/ta...
             http://www.flickr.com/photos/britishlibrary/ta...
 A., A. A. http://www.flickr.com/photos/britishlibrary/ta...
3 A., E. S. http://www.flickr.com/photos/britishlibrary/ta...
4 A., E. S. http://www.flickr.com/photos/britishlibrary/ta...
```

Alternatively, we could also remove the columns by passing them to the columns parameter directly instead of separately specifying the labels to be removed and the axis where pandas should look for the labels:

```
Python
>>> df.drop(columns=to_drop, inplace=True)
```

This syntax is more intuitive and readable. What we're trying to do here is directly apparent.

If you know in advance which columns you'd like to retain, another option is to pass them to the usecols argument of pd.read_csv.



Your Guided Tour Through the Python 3.9 Interpreter »

Remove ads

Changing the Index of a DataFrame

A pandas Index extends the functionality of NumPy arrays to allow for more versatile slicing and labeling. In many cases, it is helpful to use a uniquely valued identifying field of the data as its index.

For example, in the dataset used in the previous section, it can be expected that when a librarian searches for a record, they may input the unique identifier (values in the Identifier column) for a book:

```
Python
>>> df['Identifier'].is_unique
True
```

Let's replace the existing index with this column using set_index:

```
Python
                                                                           >>>
>>> df = df.set_index('Identifier')
>>> df.head()
                Place of Publication Date of Publication \
                                             1879 [1878]
206
                              London
            London; Virtue & Yorston
216
                                                    1868
218
                              London
                                                    1869
472
                              London
                                                    1851
480
                              London
                                                    1857
                        Publisher \
                 S. Tinsley & Co.
206
                     Virtue & Co.
216
            Bradbury, Evans & Co.
218
472
                    James Darling
             Wertheim & Macintosh
480
                                                        Title
                                                                   Author \
206
                            Walter Forbes. [A novel.] By A. A
            All for Greed. [A novel. The dedication signed... A., A. A.
216
            Love the Avenger. By the author of "All for Gr... A., A. A.
218
            Welsh Sketches, chiefly ecclesiastical, to the... A., E. S.
472
            [The World in which I live, and my place in it... A., E. S.
480
                                                   Flickr URL
            http://www.flickr.com/photos/britishlibrary/ta...
206
            http://www.flickr.com/photos/britishlibrary/ta...
216
218
            http://www.flickr.com/photos/britishlibrary/ta...
            http://www.flickr.com/photos/britishlibrary/ta...
472
            http://www.flickr.com/photos/britishlibrary/ta...
480
```

Technical Detail: Unlike primary keys in SQL, a pandas Index doesn't make any guarantee of being unique, although many indexing and merging operations will notice

We can access each record in a straightforward way with <code>loc[]</code>. Although <code>loc[]</code> may not have all that intuitive of a name, it allows us to do <code>label-based indexing</code>, which is the labeling of a row or record without regard to its position:

In other words, 206 is the first label of the index. To access it by *position*, we could use df.iloc[0], which does position-based indexing.

Technical Detail: .loc[] is technically a class instance and has some special syntax that doesn't conform exactly to most plain-vanilla Python instance methods.

Previously, our index was a RangeIndex: integers starting from 0, analogous to Python's built-in range. By passing a column name to set_index, we have changed the index to the values in Identifier.

You may have noticed that we reassigned the variable to the object returned by the method with df = df.set_index(...). This is because, by default, the method returns a modified copy of our object and does not make the changes directly to the object. We can avoid this by setting the inplace parameter:

```
df.set_index('Identifier', inplace=True)
```

Tidying up Fields in the Data

So far, we have removed unnecessary columns and changed the index of our DataFrame to something more sensible. In this section, we will clean specific columns and get them to a uniform format to get a better understanding of the dataset and enforce consistency. In particular, we will be cleaning Date of Publication and Place of Publication.

Upon inspection, all of the data types are currently the object dtype, which is roughly analogous to str in native Python.

It encapsulates any field that can't be neatly fit as numerical or categorical data. This makes sense since we're working with data that is initially a bunch of messy strings:

```
Python
>>> df.get_dtype_counts()
object 6
```

One field where it makes sense to enforce a numeric value is the date of publication so that we can do calculations down the road:

```
Python
                                                                              >>>
>>> df.loc[1905:, 'Date of Publication'].head(10)
Identifier
1905
               1888
1929
        1839, 38-54
            [1897?]
2836
2854
               1865
            1860-63
2956
2957
               1873
3017
               1866
3131
               1899
4598
               1814
               1820
4884
Name: Date of Publication, dtype: object
```

A particular book can have only one date of publication. Therefore, we need to do the following:

- Remove the extra dates in square brackets, wherever present: 1879 [1878]
- Convert date ranges to their "start date", wherever present: 1860-63; 1839, 38-54
- Completely remove the dates we are not certain about and replace them with NumPy's NaN: [1897?]
- Convert the string nan to NumPy's NaN value

Synthesizing these patterns, we can actually take advantage of a single regular expression to extract the publication year:

```
Python >>> regex = r'^(\d{4})'
```

The regular expression above is meant to find any four digits at the haginning of a string

which suffices for our case. The above is a *raw string* (meaning that a backslash is no longer an escape character), which is standard practice with regular expressions.

The \d represents any digit, and {4} repeats this rule four times. The ^ character matches the start of a string, and the parentheses denote a capturing group, which signals to pandas that we want to extract that part of the regex. (We want ^ to avoid cases where [starts off the string.)

Let's see what happens when we run this regex across our dataset:

```
Python
                                                                            >>>
>>> extr = df['Date of Publication'].str.extract(r'^(\d{4})', expand=False)
>>> extr.head()
Identifier
206
       1879
      1868
216
      1869
218
472
      1851
480
      1857
Name: Date of Publication, dtype: object
```

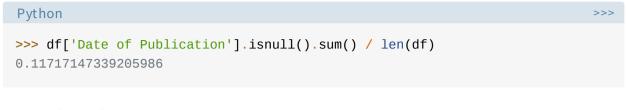
Further Reading: Not familiar with regex? You can inspect the expression above at regex101.com and learn all about regular expressions with Regular Expressions: Regexes in Python.

Technically, this column still has object dtype, but we can easily get its numerical version with pd.to_numeric:

```
Python

>>> df['Date of Publication'] = pd.to_numeric(extr)
>>> df['Date of Publication'].dtype
dtype('float64')
```

This results in about one in every ten values being missing, which is a small price to pay for now being able to do computations on the remaining valid values:



Great! That's done!



Combining str Methods with NumPy to Clean Columns

Above, you may have noticed the use of df['Date of Publication'].str. This attribute is a way to access speedy string operations in pandas that largely mimic operations on native Python strings or compiled regular expressions, such as .split(), .replace(), and .capitalize().

To clean the Place of Publication field, we can combine pandas str methods with NumPy's np.where function, which is basically a vectorized form of Excel's IF() macro. It has the following syntax:

```
Python
>>> np.where(condition, then, else)
```

Here, condition is either an array-like object or a Boolean mask. then is the value to be used if condition evaluates to True, and else is the value to be used otherwise.

Essentially, .where() takes each element in the object used for condition, checks whether that particular element evaluates to True in the context of the condition, and returns an ndarray containing then or else, depending on which applies.

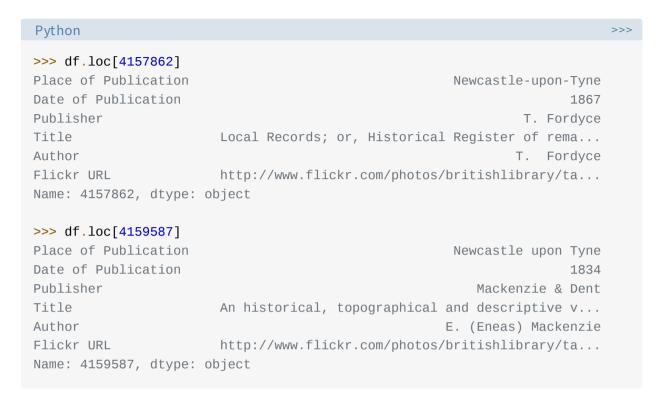
It can be nested into a compound if-then statement, allowing us to compute values based on multiple conditions:

We'll be making use of these two functions to clean Place of Publication since this column has string objects. Here are the contents of the column:

```
Python
                                                                              >>>
>>> df['Place of Publication'].head(10)
Identifier
206
                                       London
216
                    London; Virtue & Yorston
218
                                       London
472
                                       London
480
                                       London
481
                                       London
519
                                       London
667
        pp. 40. G. Bryan & Co: Oxford, 1898
874
                                      London]
1143
                                       London
Name: Place of Publication, dtype: object
```

We see that for some rows, the place of publication is surrounded by other unnecessary information. If we were to look at more values, we would see that this is the case for only some rows that have their place of publication as 'London' or 'Oxford'.

Let's take a look at two specific entries:



These two books were published in the same place, but one has hyphens in the name of the place while the other does not.

To clean this column in one sweep, we can use str.contains() to get a Boolean mask.

We clean the column as follows:

```
Python
                                                                            >>>
>>> pub = df['Place of Publication']
>>> london = pub.str.contains('London')
>>> london[:5]
Identifier
206
       True
216
       True
218
       True
472
       True
480
       True
Name: Place of Publication, dtype: bool
>>> oxford = pub.str.contains('0xford')
```

We combine them with np.where:

```
Python
                                                                            >>>
df['Place of Publication'] = np.where(london, 'London',
                                       np.where(oxford, '0xford',
                                                pub.str.replace('-', ' ')))
>>> df['Place of Publication'].head()
Identifier
206
       London
216
       London
218
       London
472
       London
480
       London
Name: Place of Publication, dtype: object
```

Here, the np.where function is called in a nested structure, with condition being a Series of Booleans obtained with str.contains(). The contains() method works similarly to the built in in knowledged to find the occurrence of an antity in an iterable (or substring in a

built-in the keyword used to find the occurrence of an entity in affice able (of substring in a string).

The replacement to be used is a string representing our desired place of publication. We also replace hyphens with a space with str.replace() and reassign to the column in our DataFrame.

Although there is more dirty data in this dataset, we will discuss only these two columns for now.

Let's have a look at the first five entries, which look a lot crisper than when we started out:

Python					>>>
>>> df.hea	ad()				
	Place of Publication Date of Pub	lication		Publishe	er \
206	London	1879	S	Tinsley & (Co.
216	London	1868	\	Virtue & Co	ο.
218	London	1869	Bradbury,	Evans & Co	Ο.
472	London	1851	Já	ames Darlir	ng
480	London	1857	Wertheim	& Macintos	sh
			Title	710101	\
206	Walter Forbes.			AA	
216	All for Greed. [A novel. The de	dication	signed	A. A A.	
218	Love the Avenger. By the author	of "All	for Gr	Α. Α Α.	
472	Welsh Sketches, chiefly ecclesi	astical,	to the	E. S A.	
480	[The World in which I live, and	my place	e in it	E. S A.	
		F	lickr URL		
206	http://www.flickr.com/photos/br	itishlibr	arv/ta		
216	http://www.flickr.com/photos/br		-		
218	http://www.flickr.com/photos/br		-		
472	http://www.flickr.com/photos/br		-		
480	http://www.flickr.com/photos/br		-		

Note: At this point, Place of Publication would be a good candidate for conversion

eterne cho ponic; i tado or i abtituactor trodia de a goda candidace roi contendior

to a Categorical dtype, because we can encode the fairly small unique set of cities with integers. (*The memory usage of a Categorical is proportional to the number of*

categories plus the length of the data; an object dtype is a constant times the length of the data.)



Online Python Training for Teams »

1 Remove ads

Cleaning the Entire Dataset Using the applymap Function

In certain situations, you will see that the "dirt" is not localized to one column but is more spread out.

There are some instances where it would be helpful to apply a customized function to each cell or element of a DataFrame. pandas .applymap() method is similar to the in-built map() function and simply applies a function to all the elements in a DataFrame.

Let's look at an example. We will create a DataFrame out of the "university_towns.txt" file:

Shell

\$ head Datasets/univerisity_towns.txt Alabama[edit] Auburn (Auburn University)[1] Florence (University of North Alabama) Jacksonville (Jacksonville State University)[2] Livingston (University of West Alabama)[2] Montevallo (University of Montevallo)[2] Troy (Troy University)[2] Tuscaloosa (University of Alabama, Stillman College, Shelton State)[3][4]

```
Tuskegee (Tuskegee University)[5]
Alaska[edit]
```

We see that we have periodic state names followed by the university towns in that state: StateA TownA1 TownA2 StateB TownB1 TownB2.... If we look at the way state names are written in the file, we'll see that all of them have the "[edit]" substring in them.

We can take advantage of this pattern by creating a *list of (state, city) tuples* and wrapping that list in a DataFrame:

```
Python
                                                                           >>>
>>> university_towns = []
>>> with open('Datasets/university_towns.txt') as file:
        for line in file:
           if '[edit]' in line:
                # Remember this `state` until the next is found
                state = line
            else:
                # Otherwise, we have a city; keep `state` as last-seen
                university_towns.append((state, line))
>>> university_towns[:5]
[('Alabama[edit]\n', 'Auburn (Auburn University)[1]\n'),
('Alabama[edit]\n', 'Florence (University of North Alabama)\n'),
 ('Alabama[edit]\n', 'Jacksonville (Jacksonville State University)[2]\n'),
('Alabama[edit]\n', 'Livingston (University of West Alabama)[2]\n'),
 ('Alabama[edit]\n', 'Montevallo (University of Montevallo)[2]\n')]
```

We can wrap this list in a DataFrame and set the columns as "State" and "RegionName". pandas will take each element in the list and set State to the left value and RegionName to the right value.

The resulting DataFrame looks like this:

```
Python
                                                                          >>>
>>> towns_df = pd.DataFrame(university_towns,
                           columns=['State', 'RegionName'])
>>> towns_df.head()
State
                                               RegionName
0 Alabama[edit]\n
                                     Auburn (Auburn University)[1]\n
1 Alabama[edit]\n
                            Florence (University of North Alabama)\n
 Alabama[edit]\n Jacksonville (Jacksonville State University)[2]\n
3 Alabama[edit]\n
                        Livingston (University of West Alabama)[2]\n
4 Alabama[edit]\n
                          Montevallo (University of Montevallo)[2]\n
```

While we could have cleaned these strings in the for loop above, pandas makes it easy. We only need the state name and the town name and can remove everything else. While we could use pandas'.str() methods again here, we could also use applymap() to map a Python callable to each element of the DataFrame.

We have been using the term *element*, but what exactly do we mean by it? Consider the following "toy" DataFrame:

```
Python >>>

0 1
0 Mock Dataset
1 Python pandas
2 Real Python
3 NumPy Clean
```

In this example, each cell ('Mock', 'Dataset', 'Python', 'pandas', etc.) is an element. Therefore, applymap() will apply a function to each of these independently. Let's define that function:

```
Python

>>> def get_citystate(item):
...    if ' (' in item:
...        return item[:item.find(' (')]
...    elif '[' in item:
...        return item[:item.find('[')]
...    else:
...    return item
```

pandas' .applymap() only takes one parameter, which is the function (callable) that should be applied to each element:

```
Python
>>> towns_df = towns_df.applymap(get_citystate)
```

First, we define a Python function that takes an element from the DataFrame as its parameter. Inside the function, checks are performed to determine whether there's a (or [in the element or not.

Depending on the check, values are returned accordingly by the function. Finally, the applymap() function is called on our object. Now the DataFrame is much neater:

```
Python

>>> towns_df.head()
    State RegionName
0 Alabama Auburn
1 Alabama Florence
2 Alabama Jacksonville
3 Alabama Livingston
4 Alabama Montevallo
```

The applymap() method took each element from the DataFrame, passed it to the function, and the original value was replaced by the returned value. It's that simple!

Technical Detail: While it is a convenient and versatile method, .applymap can have significant runtime for larger datasets, because it maps a Python callable to each individual element. In some cases, it can be more efficient to do *vectorized* operations that utilize Cython or NumPY (which, in turn, makes calls in C) under the hood.



Remove ads

Renaming Columns and Skipping Rows

Often, the datasets you'll work with will have either column names that are not easy to understand, or unimportant information in the first few and/or last rows, such as definitions of the terms in the dataset, or footnotes.

In that case, we'd want to rename columns and skip certain rows so that we can drill down to necessary information with correct and sensible labels.

To demonstrate how we can go about doing this, let's first take a glance at the initial five rows of the "olympics.csv" dataset:

\$ head -n 5 Datasets/olympics.csv 0,1,2,3,4,5,6,7,8,9,10,11,12,13,14,15 ,? Summer,01 !,02 !,03 !,Total,? Winter,01 !,02 !,03 !,Total,? Games,01 !,02 ! Afghanistan (AFG),13,0,0,2,2,0,0,0,0,13,0,0,2,2 Algeria (ALG),12,5,2,8,15,3,0,0,0,0,15,5,2,8,15 Argentina (ARG),23,18,24,28,70,18,0,0,0,0,41,18,24,28,70

Now, we'll read it into a pandas DataFrame:

```
Python
                                                                            >>>
>>> olympics_df = pd.read_csv('Datasets/olympics.csv')
>>> olympics_df.head()
                                                                              8
0
                     ? Summer
                                01 !
                                       02!
                                             03 !
                                                  Total
                                                          ? Winter
                                                       2
  Afghanistan (AFG)
                            13
                                                                              0
                            12
       Algeria (ALG)
                                                8
                                                      15
                                                                  3
                                                                              0
    Argentina (ARG)
                             23
                                   18
                                         24
                                                      70
                                                                 18
3
                                               28
                                                                              0
      Armenia (ARM)
                                          2
                                   1
                                                9
                                                      12
                                                                  6
                                                                        0
                                                                              0
            10
                                  13
                                                        15
      9
                     11
                           12
                                        14
                         01 !
                               02!
                                      03!
                                            Combined total
               ? Games
                                         2
                     13
                                   0
                                         8
                                                        15
                     15
                                   2
      0
             0
                     41
                           18
                                  24
                                        28
                                                        70
                            1
                                   2
      0
                     11
                                         9
                                                        12
```

This is messy indeed! The columns are the string form of integers indexed at 0. The row which should have been our header (i.e. the one to be used to set the column names) is at $olympics_df.iloc[0]$. This happened because our CSV file starts with 0, 1, 2, ..., 15.

Also, if we were to go to the source of this dataset, we'd see that NaN above should really be something like "Country",? Summer is supposed to represent "Summer Games", 01! should be "Gold", and so on.

Therefore, we need to do two things:

- Skip one row and set the header as the first (0-indexed) row
- Rename the columns

We can skip rows and set the header while reading the CSV file by passing some parameters to the read_csv() function.

This function takes *a lot* of optional parameters, but in this case we only need one (header) to remove the 0th row:

Ру	thon										>>>
>>> olympics_df = pd.read_csv('Datasets/olympics.csv', header=1)											
>>> olympics_df.head()											
				nmer 01 !		03 !	Total		inter \		
0	А	_	an (AFG)			0	2	2		0	
1		•	ia (ALG)			2	8	15	5	3	
2		_	na (ARG)			24	28	70	9	18	
3			ia (ARM)			2	9	12	2	6	
4	Austral	asia (AN	Z) [ANZ]	2	3	4	5	12	2	0	
	04 1 4	00 1 4	00 1 4	T.1.1.4	0.00	0.1				,	
	01 !.1			Total.1						\	
0	0	0	0	0	13		0	0	2		
1	0	0	0	0	15		5	2	8		
2	0	0	0	0	41		18	24	28		
3	0	0	0	0	11		1	2	9		
4	0	0	0	0	2		3	4	5		
Combined total											
0		2									
1		15									
2		70									
3		12									
4		12									

We now have the correct row set as the header and all unnecessary rows removed. Take note of how pandas has changed the name of the column containing the name of the countries from NaN to Unnamed: 0.

To rename the columns, we will make use of a DataFrame's rename() method, which allows you to relabel an axis based on a *mapping* (in this case, a dict).

Let's start by defining a dictionary that maps current column names (as keys) to more usable ones (the dictionary's values):

We call the rename() function on our object:

```
Python
>>> olympics_df.rename(columns=new_names, inplace=True)
```

Setting *inplace* to True specifies that our changes be made directly to the object. Let's see if this checks out:

Ру	rthon .									>>>	
>>> olympics_df.head()											
		(Country	Summer 0]	ympics	Gold	Silver	Bronze	Total	\	
0	Afg	hanistar	n (AFG)		13	Θ	0	2	2		
1		Algeria	a (ALG)		12	5	2	8	15		
2	А	rgentina	a (ARG)		23	18	24	28	70		
3		Armenia	a (ARM)		5	1	2	9	12		
4	Australas	ia (ANZ)) [ANZ]		2	3	4	5	12		
	Winter Ol	ympics	Gold.1	Silver.1	Bronze	e.1 To	otal.1 #	Games	Gold.2	\	
0		0	0	0		Θ	0	13	Θ		
1		3	0	0		0	0	15	5		
2		18	0	0		0	0	41	18		
3		6	0	0		0	0	11	1		
4		Θ	0	0		0	0	2	3		
	Silver.2	Bronze	.2 Comb	ined total							
0	0	5. 311201	2	2							
1	2		8	15							
2	24	2	28	70							
3	2		9	12							
4	4		5	12	-						

Find Your Dream Python Job

pythoniobsha.com



Remove ads

Python Data Cleaning: Recap and Resources

In this tutorial, you learned how you can drop unnecessary information from a dataset using the drop() function, as well as how to set an index for your dataset so that items in it can

be referenced easily.

Moreover, you learned how to clean object fields with the .str() accessor and how to clean the entire dataset using the applymap() method. Lastly, we explored how to skip rows in a CSV file and rename columns using the rename() method.

Knowing about data cleaning is very important, because it is a big part of data science. You now have a basic understanding of how pandas and NumPy can be leveraged to clean datasets!

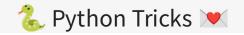
Check out the links below to find additional resources that will help you on your Python data science journey:

- The pandas documentation
- The NumPy documentation
- Python for Data Analysis by Wes McKinney, the creator of pandas
- pandas Cookbook by Ted Petrou, a data science trainer and consultant

Free Bonus: Click here to get access to a free NumPy Resources Guide that points you to the best tutorials, videos, and books for improving your NumPy skills.



Watch Now This tutorial has a related video course created by the Real Python team. Watch it together with the written tutorial to deepen your understanding: **Data**Cleaning With pandas and NumPy



Get a short & sweet **Python Trick** delivered to your inbox every couple of days. No spam

ever. Unsubscribe any time. Curated by the Real Python team.

```
1# How to merge two dicts
2# in Python 3.5+
3
4>>> x = {'a': 1, 'b': 2}
5>>> y = {'b': 3, 'c': 4}
6
7>>> z = {**x, **y}
8
9>>> z
10 {'c': 4, 'a': 1, 'b': 3}
```

Email Address

Send Me Python Tricks »

About Malay Agarwal

A tech geek with a philosophical mind and a hand that can wield a pen.

» More about Malay

Each tutorial at Real Python is created by a team of developers so that it meets our high quality standards. The team members who worked on this tutorial are:

Master <u>Real-World Python Skills</u> With Unlimited Access to Real Python

Join us and get access to thousands of tutorials, hands-on video courses, and a community of expert Pythonistas:

Level Up Your Python Skills »

What Do You Think?

Rate this article:







What's your #1 takeaway or favorite thing you learned? How are you going to put your newfound skills to use? Leave a comment below and let us know.

Commenting Tips: The most useful comments are those written with the goal of learning from or helping out other students. Get tips for asking good questions and get answers to common questions in our support portal.

Looking for a real-time conversation? Visit the Real Python Community Chat or join the next "Office Hours" Live Q&A Session. Happy Pythoning!

Keep Learning

Related Tutorial Categories: data-science intermediate

Recommended Video Course: Data Cleaning With pandas and NumPy



Your Practical Introduction to Python 3 »

Remove ads

© 2012–2023 Real Python \cdot Newsletter \cdot Podcast \cdot YouTube \cdot Twitter \cdot Facebook \cdot Instagram \cdot Python Tutorials \cdot Search \cdot Privacy Policy \cdot Energy Policy \cdot Advertise \cdot Contact

