

bonus_top_25_tricks

February 29, 2024

1 Bonus video: My top 25 pandas tricks

Full course: [pandas in 30 days](#)

© 2024 Data School. All rights reserved.

1. Show installed versions
2. Create an example DataFrame
3. Rename columns
4. Reverse row order
5. Reverse column order
6. Select columns by data type
7. Convert strings to numbers
8. Reduce DataFrame size
9. Build a DataFrame from multiple files (row-wise)
10. Build a DataFrame from multiple files (column-wise)
11. Create a DataFrame from the clipboard
12. Split a DataFrame into two random subsets
13. Filter a DataFrame by multiple categories
14. Filter a DataFrame by largest categories
15. Handle missing values
16. Split a string into multiple columns
17. Expand a Series of lists into a DataFrame
18. Aggregate by multiple functions
19. Combine the output of an aggregation with a DataFrame
20. Select a slice of rows and columns
21. Reshape a MultiIndexed Series
22. Create a pivot table
23. Convert continuous data into categorical data
24. Change display options
25. Style a DataFrame
26. Bonus trick: Profile a DataFrame

1.1 Load example datasets

```
[1]: import pandas as pd
import numpy as np
```

```
[2]: drinks = pd.read_csv('http://bit.ly/drinksbycountry')
     movies = pd.read_csv('http://bit.ly/imdbratings')
     orders = pd.read_csv('http://bit.ly/chiporders', sep='\t')
     orders['item_price'] = orders.item_price.str.replace('$', '').astype('float')
     stocks = pd.read_csv('http://bit.ly/smallstocks', parse_dates=['Date'])
     titanic = pd.read_csv('http://bit.ly/kaggletrain')
     ufo = pd.read_csv('http://bit.ly/uforeports', parse_dates=['Time'])
```

1.2 1. Show installed versions

Sometimes you need to know the pandas version you're using, especially when reading the pandas documentation. You can show the pandas version by typing:

```
[3]: pd.__version__
```

```
[3]: '2.1.4'
```

But if you also need to know the versions of pandas' dependencies, you can use the `show_versions()` function:

```
[4]: pd.show_versions()
```

```
/Users/kevin/miniconda3/envs/2312/lib/python3.11/site-
packages/_distutils_hack/__init__.py:33: UserWarning: Setuptools is replacing
distutils.
  warnings.warn("Setuptools is replacing distutils.")
```

INSTALLED VERSIONS

```
-----
commit           : a671b5a8bf5dd13fb19f0e88edc679bc9e15c673
python           : 3.11.5.final.0
python-bits      : 64
OS               : Darwin
OS-release       : 22.5.0
Version          : Darwin Kernel Version 22.5.0: Mon Apr 24 20:51:50 PDT
2023; root:xnu-8796.121.2~5/RELEASE_X86_64
machine          : x86_64
processor        : i386
byteorder        : little
LC_ALL           : None
LANG             : en_US.UTF-8
LOCALE           : en_US.UTF-8

pandas           : 2.1.4
numpy            : 1.24.3
pytz             : 2023.3.post1
dateutil         : 2.8.2
setuptools       : 68.2.2
```

pip	: 23.3.1
Cython	: None
pytest	: None
hypothesis	: None
sphinx	: None
blosc	: None
feather	: None
xlsxwriter	: None
lxml.etree	: None
html5lib	: None
pymysql	: None
psycpg2	: None
jinja2	: 3.1.2
IPython	: 8.15.0
pandas_datareader	: None
bs4	: 4.12.2
bottleneck	: 1.3.5
dataframe-api-compat	: None
fastparquet	: None
fsspec	: None
gcsfs	: None
matplotlib	: 3.8.0
numba	: 0.58.1
numexpr	: 2.8.7
odfpy	: None
openpyxl	: None
pandas_gbq	: None
pyarrow	: None
pyreadstat	: None
pyxlsb	: None
s3fs	: None
scipy	: 1.11.4
sqlalchemy	: None
tables	: None
tabulate	: None
xarray	: None
xlrd	: None
zstandard	: None
tzdata	: 2023.3
qtpy	: 2.4.1
pyqt5	: None

You can see the versions of Python, pandas, NumPy, matplotlib, and more.

1.3 2. Create an example DataFrame

Let's say that you want to demonstrate some pandas code. You need an example DataFrame to work with.

There are many ways to do this, but my favorite way is to pass a dictionary to the `DataFrame()` constructor, in which the dictionary keys are the column names and the dictionary values are lists of column values:

```
[5]: df = pd.DataFrame({'col one':[100, 200], 'col two':[300, 400]})
df
```

```
[5]:   col one  col two
0      100     300
1      200     400
```

Now if you need a much larger DataFrame, the above method will require way too much typing. In that case, you can use NumPy's `random.rand()` function, tell it the number of rows and columns, and pass that to the `DataFrame()` constructor:

```
[6]: pd.DataFrame(np.random.rand(4, 8))
```

```
[6]:      0      1      2      3      4      5      6  \
0  0.097242  0.866827  0.313935  0.098417  0.737658  0.375012  0.800957
1  0.789413  0.712345  0.594955  0.297288  0.187803  0.273729  0.302122
2  0.553280  0.345825  0.880758  0.855058  0.556704  0.615995  0.846177
3  0.685043  0.453427  0.965297  0.869912  0.275472  0.082029  0.823650

      7
0  0.454138
1  0.338817
2  0.644631
3  0.533750
```

That's pretty good, but if you also want non-numeric column names, you can coerce a string of letters to a list and then pass that list to the `columns` parameter:

```
[7]: pd.DataFrame(np.random.rand(4, 8), columns=list('abcdefgh'))
```

```
[7]:      a      b      c      d      e      f      g  \
0  0.094968  0.378716  0.080192  0.333718  0.905132  0.161064  0.709727
1  0.854602  0.722020  0.410714  0.561996  0.504466  0.593359  0.107510
2  0.992585  0.777179  0.872854  0.892928  0.465297  0.778264  0.351882
3  0.674122  0.269331  0.988019  0.506475  0.823134  0.839451  0.722780

      h
0  0.503617
1  0.754404
2  0.866033
3  0.105100
```

As you might guess, your string will need to have the same number of characters as there are columns.

1.4 3. Rename columns

Let's take a look at the example DataFrame we created in the last trick:

```
[8]: df
```

```
[8]:   col one  col two
0      100     300
1      200     400
```

I prefer to use dot notation to select pandas columns, but that won't work since the column names have spaces. Let's fix this.

The most flexible method for renaming columns is the `rename()` method. You pass it a dictionary in which the keys are the old names and the values are the new names, and you also specify the axis:

```
[9]: df = df.rename({'col one': 'col_one', 'col two': 'col_two'}, axis='columns')
```

The best thing about this method is that you can use it to rename any number of columns, whether it be just one column or all columns.

Now if you're going to rename all of the columns at once, a simpler method is just to overwrite the columns attribute of the DataFrame:

```
[10]: df.columns = ['col_one', 'col_two']
```

Now if the only thing you're doing is replacing spaces with underscores, an even better method is to use the `str.replace()` method, since you don't have to type out all of the column names:

```
[11]: df.columns = df.columns.str.replace(' ', '_')
```

All three of these methods have the same result, which is to rename the columns so that they don't have any spaces:

```
[12]: df
```

```
[12]:   col_one  col_two
0      100     300
1      200     400
```

Finally, if you just need to add a prefix or suffix to all of your column names, you can use the `add_prefix()` method...

```
[13]: df.add_prefix('X_')
```

```
[13]:   X_col_one  X_col_two
0         100         300
1         200         400
```

...or the `add_suffix()` method:

```
[14]: df.add_suffix('_Y')
```

```
[14]:   col_one_Y  col_two_Y
      0      100      300
      1      200      400
```

1.5 4. Reverse row order

Let's take a look at the drinks DataFrame:

```
[15]: drinks.head()
```

```
[15]:   country  beer_servings  spirit_servings  wine_servings  \
0  Afghanistan           0              0              0
1   Albania           89             132             54
2   Algeria           25              0             14
3   Andorra          245             138            312
4    Angola          217              57             45

   total_litres_of_pure_alcohol  continent
0                        0.0      Asia
1                        4.9    Europe
2                        0.7    Africa
3                       12.4    Europe
4                        5.9    Africa
```

This is a dataset of average alcohol consumption by country. What if you wanted to reverse the order of the rows?

The most straightforward method is to use the `loc` accessor and pass it `::-1`, which is the same slicing notation used to reverse a Python list:

```
[16]: drinks.loc[::-1].head()
```

```
[16]:   country  beer_servings  spirit_servings  wine_servings  \
192  Zimbabwe           64              18              4
191   Zambia           32              19              4
190   Yemen            6              0              0
189  Vietnam          111              2              1
188  Venezuela          333             100              3

   total_litres_of_pure_alcohol  continent
192                        4.7      Africa
191                        2.5      Africa
190                        0.1      Asia
189                        2.0      Asia
188                        7.7  South America
```

What if you also wanted to reset the index so that it starts at zero?

You would use the `reset_index()` method and tell it to drop the old index entirely:

```
[17]: drinks.loc[::-1].reset_index(drop=True).head()
```

```
[17]:
```

	country	beer_servings	spirit_servings	wine_servings	\
0	Zimbabwe	64	18	4	
1	Zambia	32	19	4	
2	Yemen	6	0	0	
3	Vietnam	111	2	1	
4	Venezuela	333	100	3	

	total_litres_of_pure_alcohol	continent
0	4.7	Africa
1	2.5	Africa
2	0.1	Asia
3	2.0	Asia
4	7.7	South America

As you can see, the rows are in reverse order but the index has been reset to the default integer index.

1.6 5. Reverse column order

Similar to the previous trick, you can also use `loc` to reverse the left-to-right order of your columns:

```
[18]: drinks.loc[:, ::-1].head()
```

```
[18]:
```

	continent	total_litres_of_pure_alcohol	wine_servings	spirit_servings	\
0	Asia	0.0	0	0	
1	Europe	4.9	54	132	
2	Africa	0.7	14	0	
3	Europe	12.4	312	138	
4	Africa	5.9	45	57	

	beer_servings	country
0	0	Afghanistan
1	89	Albania
2	25	Algeria
3	245	Andorra
4	217	Angola

The colon before the comma means “select all rows”, and the `::-1` after the comma means “reverse the columns”, which is why “country” is now on the right side.

1.7 6. Select columns by data type

Here are the data types of the drinks DataFrame:

```
[19]: drinks.dtypes
```

```
[19]: country          object
      beer_servings    int64
      spirit_servings   int64
      wine_servings     int64
      total_litres_of_pure_alcohol float64
      continent         object
      dtype: object
```

Let's say you need to select only the numeric columns. You can use the `select_dtypes()` method:

```
[20]: drinks.select_dtypes(include='number').head()
```

```
[20]:   beer_servings  spirit_servings  wine_servings  total_litres_of_pure_alcohol
0             0             0             0             0.0
1            89            132            54             4.9
2            25             0            14             0.7
3           245            138           312            12.4
4           217             57            45             5.9
```

This includes both int and float columns.

You could also use this method to select just the object columns:

```
[21]: drinks.select_dtypes(include='object').head()
```

```
[21]:   country continent
0  Afghanistan      Asia
1    Albania     Europe
2    Algeria     Africa
3   Andorra     Europe
4    Angola     Africa
```

You can tell it to include multiple data types by passing a list:

```
[22]: drinks.select_dtypes(include=['number', 'object', 'category', 'datetime']).
      ↪head()
```

```
[22]:   country  beer_servings  spirit_servings  wine_servings  \
0  Afghanistan             0             0             0
1    Albania             89            132            54
2    Algeria             25             0            14
3   Andorra            245            138           312
4    Angola            217             57            45

      total_litres_of_pure_alcohol  continent
0                  0.0      Asia
1                  4.9     Europe
2                  0.7     Africa
3                 12.4     Europe
4                  5.9     Africa
```


You can also tell it to exclude certain data types:

```
[23]: drinks.select_dtypes(exclude='number').head()
```

```
[23]:      country continent
0  Afghanistan      Asia
1     Albania    Europe
2     Algeria    Africa
3     Andorra    Europe
4      Angola    Africa
```

1.8 7. Convert strings to numbers

Let's create another example DataFrame:

```
[24]: df = pd.DataFrame({'col_one':['1.1', '2.2', '3.3'],
                        'col_two':['4.4', '5.5', '6.6'],
                        'col_three':['7.7', '8.8', '-']})
df
```

```
[24]:   col_one col_two col_three
0      1.1     4.4      7.7
1      2.2     5.5      8.8
2      3.3     6.6      -
```

These numbers are actually stored as strings, which results in object columns:

```
[25]: df.dtypes
```

```
[25]: col_one      object
col_two      object
col_three    object
dtype: object
```

In order to do mathematical operations on these columns, we need to convert the data types to numeric. You can use the `astype()` method on the first two columns:

```
[26]: df.astype({'col_one':'float', 'col_two':'float'}).dtypes
```

```
[26]: col_one      float64
col_two      float64
col_three    object
dtype: object
```

However, this would have resulted in an error if you tried to use it on the third column, because that column contains a dash to represent zero and pandas doesn't understand how to handle it.

Instead, you can use the `to_numeric()` function on the third column and tell it to convert any invalid input into NaN values:

```
[27]: pd.to_numeric(df.col_three, errors='coerce')
```

```
[27]: 0    7.7
      1    8.8
      2   NaN
      Name: col_three, dtype: float64
```

If you know that the NaN values actually represent zeros, you can fill them with zeros using the `fillna()` method:

```
[28]: pd.to_numeric(df.col_three, errors='coerce').fillna(0)
```

```
[28]: 0    7.7
      1    8.8
      2    0.0
      Name: col_three, dtype: float64
```

Finally, you can apply this function to the entire DataFrame all at once by using the `apply()` method:

```
[29]: df = df.apply(pd.to_numeric, errors='coerce').fillna(0)
      df
```

```
[29]:   col_one  col_two  col_three
0      1.1      4.4        7.7
1      2.2      5.5        8.8
2      3.3      6.6        0.0
```

This one line of code accomplishes our goal, because all of the data types have now been converted to float:

```
[30]: df.dtypes
```

```
[30]: col_one      float64
      col_two      float64
      col_three      float64
      dtype: object
```

1.9 8. Reduce DataFrame size

pandas DataFrames are designed to fit into memory, and so sometimes you need to reduce the DataFrame size in order to work with it on your system.

Here's the size of the drinks DataFrame:

```
[31]: drinks.info(memory_usage='deep')
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 193 entries, 0 to 192
Data columns (total 6 columns):
#   Column                                Non-Null Count  Dtype
#   ...
```

```

---  -----
0   country                193 non-null    object
1   beer_servings          193 non-null    int64
2   spirit_servings         193 non-null    int64
3   wine_servings          193 non-null    int64
4   total_litres_of_pure_alcohol  193 non-null    float64
5   continent              193 non-null    object
dtypes: float64(1), int64(3), object(2)
memory usage: 30.5 KB

```

You can see that it currently uses 30.5 KB.

If you're having performance problems with your DataFrame, or you can't even read it into memory, there are two easy steps you can take during the file reading process to reduce the DataFrame size.

The first step is to only read in the columns that you actually need, which we specify with the “usecols” parameter:

```
[32]: cols = ['beer_servings', 'continent']
small_drinks = pd.read_csv('http://bit.ly/drinksbycountry', usecols=cols)
small_drinks.info(memory_usage='deep')
```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 193 entries, 0 to 192
Data columns (total 2 columns):
#   Column          Non-Null Count  Dtype
---  ---
0   beer_servings    193 non-null    int64
1   continent        193 non-null    object
dtypes: int64(1), object(1)
memory usage: 13.7 KB

```

By only reading in these two columns, we've reduced the DataFrame size to 13.7 KB.

The second step is to convert any object columns containing categorical data to the category data type, which we specify with the “dtype” parameter:

```
[33]: dtypes = {'continent': 'category'}
smaller_drinks = pd.read_csv('http://bit.ly/drinksbycountry', usecols=cols,
                               dtype=dtypes)
smaller_drinks.info(memory_usage='deep')
```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 193 entries, 0 to 192
Data columns (total 2 columns):
#   Column          Non-Null Count  Dtype
---  ---
0   beer_servings    193 non-null    int64
1   continent        193 non-null    category
dtypes: category(1), int64(1)
memory usage: 2.4 KB

```

By reading in the continent column as the category data type, we've further reduced the DataFrame size to 2.4 KB.

Keep in mind that the category data type will only reduce memory usage if you have a small number of categories relative to the number of rows.

1.10 9. Build a DataFrame from multiple files (row-wise)

Let's say that your dataset is spread across multiple files, but you want to read the dataset into a single DataFrame.

For example, I have a small dataset of stock data in which each CSV file only includes a single day. Here's the first day:

```
[34]: pd.read_csv('data/stocks1.csv')
```

```
[34]:
```

	Date	Close	Volume	Symbol
0	2016-10-03	31.50	14070500	CSCO
1	2016-10-03	112.52	21701800	AAPL
2	2016-10-03	57.42	19189500	MSFT

Here's the second day:

```
[35]: pd.read_csv('data/stocks2.csv')
```

```
[35]:
```

	Date	Close	Volume	Symbol
0	2016-10-04	113.00	29736800	AAPL
1	2016-10-04	57.24	20085900	MSFT
2	2016-10-04	31.35	18460400	CSCO

And here's the third day:

```
[36]: pd.read_csv('data/stocks3.csv')
```

```
[36]:
```

	Date	Close	Volume	Symbol
0	2016-10-05	57.64	16726400	MSFT
1	2016-10-05	31.59	11808600	CSCO
2	2016-10-05	113.05	21453100	AAPL

You could read each CSV file into its own DataFrame, combine them together, and then delete the original DataFrames, but that would be memory inefficient and require a lot of code.

A better solution is to use the built-in glob module:

```
[37]: from glob import glob
```

You can pass a pattern to `glob()`, including wildcard characters, and it will return a list of all files that match that pattern.

In this case, glob is looking in the “data” subdirectory for all CSV files that start with the word “stocks”:

```
[38]: stock_files = sorted(glob('data/stocks?*.csv'))
      stock_files
```

```
[38]: ['data/stocks1.csv', 'data/stocks2.csv', 'data/stocks3.csv']
```

glob returns filenames in an arbitrary order, which is why we sorted the list using Python's built-in `sorted()` function.

We can then use a generator expression to read each of the files using `read_csv()` and pass the results to the `concat()` function, which will concatenate the rows into a single DataFrame:

```
[39]: pd.concat((pd.read_csv(file) for file in stock_files))
```

```
[39]:
```

	Date	Close	Volume	Symbol
0	2016-10-03	31.50	14070500	CSCO
1	2016-10-03	112.52	21701800	AAPL
2	2016-10-03	57.42	19189500	MSFT
0	2016-10-04	113.00	29736800	AAPL
1	2016-10-04	57.24	20085900	MSFT
2	2016-10-04	31.35	18460400	CSCO
0	2016-10-05	57.64	16726400	MSFT
1	2016-10-05	31.59	11808600	CSCO
2	2016-10-05	113.05	21453100	AAPL

Unfortunately, there are now duplicate values in the index. To avoid that, we can tell the `concat()` function to ignore the index and instead use the default integer index:

```
[40]: pd.concat((pd.read_csv(file) for file in stock_files), ignore_index=True)
```

```
[40]:
```

	Date	Close	Volume	Symbol
0	2016-10-03	31.50	14070500	CSCO
1	2016-10-03	112.52	21701800	AAPL
2	2016-10-03	57.42	19189500	MSFT
3	2016-10-04	113.00	29736800	AAPL
4	2016-10-04	57.24	20085900	MSFT
5	2016-10-04	31.35	18460400	CSCO
6	2016-10-05	57.64	16726400	MSFT
7	2016-10-05	31.59	11808600	CSCO
8	2016-10-05	113.05	21453100	AAPL

1.11 10. Build a DataFrame from multiple files (column-wise)

The previous trick is useful when each file contains rows from your dataset. But what if each file instead contains columns from your dataset?

Here's an example in which the drinks dataset has been split into two CSV files, and each file contains three columns:

```
[41]: pd.read_csv('data/drinks1.csv').head()
```

```
[41]:
```

	country	beer_servings	spirit_servings
0	Afghanistan	0	0
1	Albania	89	132
2	Algeria	25	0
3	Andorra	245	138
4	Angola	217	57

```
[42]: pd.read_csv('data/drinks2.csv').head()
```

```
[42]:
```

	wine_servings	total_litres_of_pure_alcohol	continent
0	0	0.0	Asia
1	54	4.9	Europe
2	14	0.7	Africa
3	312	12.4	Europe
4	45	5.9	Africa

Similar to the previous trick, we'll start by using `glob()`:

```
[43]: drink_files = sorted(glob('data/drinks?*.csv'))
```

And this time, we'll tell the `concat()` function to concatenate along the columns axis:

```
[44]: pd.concat((pd.read_csv(file) for file in drink_files), axis='columns').head()
```

```
[44]:
```

	country	beer_servings	spirit_servings	wine_servings	\
0	Afghanistan	0	0	0	
1	Albania	89	132	54	
2	Algeria	25	0	14	
3	Andorra	245	138	312	
4	Angola	217	57	45	

	total_litres_of_pure_alcohol	continent
0	0.0	Asia
1	4.9	Europe
2	0.7	Africa
3	12.4	Europe
4	5.9	Africa

Now our DataFrame has all six columns.

1.12 11. Create a DataFrame from the clipboard

Let's say that you have some data stored in an Excel spreadsheet or a [Google Sheet](#), and you want to get it into a DataFrame as quickly as possible.

Just select the data and copy it to the clipboard. Then, you can use the `read_clipboard()` function to read it into a DataFrame:

```
[45]: df = pd.read_clipboard()
df
```

```
[45]:
```

	Column A	Column B	Column C
0	1	4.4	seven
1	2	5.5	eight
2	3	6.6	nine

Just like the `read_csv()` function, `read_clipboard()` automatically detects the correct data type for each column:

```
[46]: df.dtypes
```

```
[46]: Column A      int64
      Column B    float64
      Column C     object
      dtype: object
```

Let's copy one other dataset to the clipboard:

```
[47]: df = pd.read_clipboard()
      df
```

```
[47]:
```

	Left	Right
Alice	10	40
Bob	20	50
Charlie	30	60

Amazingly, pandas has even identified the first column as the index:

```
[48]: df.index
```

```
[48]: Index(['Alice', 'Bob', 'Charlie'], dtype='object')
```

Keep in mind that if you want your work to be reproducible in the future, `read_clipboard()` is not the recommended approach.

1.13 12. Split a DataFrame into two random subsets

Let's say that you want to split a DataFrame into two parts, randomly assigning 75% of the rows to one DataFrame and the other 25% to a second DataFrame.

For example, we have a DataFrame of movie ratings with 979 rows:

```
[49]: len(movies)
```

```
[49]: 979
```

We can use the `sample()` method to randomly select 75% of the rows and assign them to the "movies_1" DataFrame:

```
[50]: movies_1 = movies.sample(frac=0.75, random_state=1234)
```

Then we can use the `drop()` method to drop all rows that are in "movies_1" and assign the remaining rows to "movies_2":

```
[51]: movies_2 = movies.drop(movies_1.index)
```

You can see that the total number of rows is correct:

```
[52]: len(movies_1) + len(movies_2)
```

```
[52]: 979
```

And you can see from the index that every movie is in either “movies_1”:

```
[53]: movies_1.index.sort_values()
```

```
[53]: Index([ 0,  2,  5,  6,  7,  8,  9, 11, 13, 16,
...
          966, 967, 969, 971, 972, 974, 975, 976, 977, 978],
          dtype='int64', length=734)
```

...or “movies_2”:

```
[54]: movies_2.index.sort_values()
```

```
[54]: Index([ 1,  3,  4, 10, 12, 14, 15, 18, 26, 30,
...
          931, 934, 937, 941, 950, 954, 960, 968, 970, 973],
          dtype='int64', length=245)
```

Keep in mind that this approach will not work if your index values are not unique.

1.14 13. Filter a DataFrame by multiple categories

Let’s take a look at the movies DataFrame:

```
[55]: movies.head()
```

```
[55]:
```

	star_rating	title	content_rating	genre	duration	\
0	9.3	The Shawshank Redemption	R	Crime	142	
1	9.2	The Godfather	R	Crime	175	
2	9.1	The Godfather: Part II	R	Crime	200	
3	9.0	The Dark Knight	PG-13	Action	152	
4	8.9	Pulp Fiction	R	Crime	154	


```
actors_list
```

0	[u'Tim Robbins', u'Morgan Freeman', u'Bob Gunt...
1	[u'Marlon Brando', u'Al Pacino', u'James Caan']
2	[u'Al Pacino', u'Robert De Niro', u'Robert Duv...
3	[u'Christian Bale', u'Heath Ledger', u'Aaron E...
4	[u'John Travolta', u'Uma Thurman', u'Samuel L...

One of the columns is genre:

```
[56]: movies.genre.unique()
```



```
[56]: array(['Crime', 'Action', 'Drama', 'Western', 'Adventure', 'Biography',
        'Comedy', 'Animation', 'Mystery', 'Horror', 'Film-Noir', 'Sci-Fi',
        'History', 'Thriller', 'Family', 'Fantasy'], dtype=object)
```

If we wanted to filter the DataFrame to only show movies with the genre Action or Drama or Western, we could use multiple conditions separated by the “or” operator:

```
[57]: movies[(movies.genre == 'Action') |
            (movies.genre == 'Drama') |
            (movies.genre == 'Western')].head()
```

```
[57]:
```

	star_rating	title	content_rating	genre	\
3	9.0	The Dark Knight	PG-13	Action	
5	8.9	12 Angry Men	NOT RATED	Drama	
6	8.9	The Good, the Bad and the Ugly	NOT RATED	Western	
9	8.9	Fight Club	R	Drama	
11	8.8	Inception	PG-13	Action	

	duration	actors_list
3	152	[u'Christian Bale', u'Heath Ledger', u'Aaron E...
5	96	[u'Henry Fonda', u'Lee J. Cobb', u'Martin Bals...
6	161	[u'Clint Eastwood', u'Eli Wallach', u'Lee Van ...
9	139	[u'Brad Pitt', u'Edward Norton', u'Helena Bonh...
11	148	[u'Leonardo DiCaprio', u'Joseph Gordon-Levitt'...

However, you can actually rewrite this code more clearly by using the `isin()` method and passing it a list of genres:

```
[58]: movies[movies.genre.isin(['Action', 'Drama', 'Western'])].head()
```

```
[58]:
```

	star_rating	title	content_rating	genre	\
3	9.0	The Dark Knight	PG-13	Action	
5	8.9	12 Angry Men	NOT RATED	Drama	
6	8.9	The Good, the Bad and the Ugly	NOT RATED	Western	
9	8.9	Fight Club	R	Drama	
11	8.8	Inception	PG-13	Action	

	duration	actors_list
3	152	[u'Christian Bale', u'Heath Ledger', u'Aaron E...
5	96	[u'Henry Fonda', u'Lee J. Cobb', u'Martin Bals...
6	161	[u'Clint Eastwood', u'Eli Wallach', u'Lee Van ...
9	139	[u'Brad Pitt', u'Edward Norton', u'Helena Bonh...
11	148	[u'Leonardo DiCaprio', u'Joseph Gordon-Levitt'...

And if you want to reverse this filter, so that you are excluding (rather than including) those three genres, you can put a tilde in front of the condition:

```
[59]: movies[~movies.genre.isin(['Action', 'Drama', 'Western'])].head()
```

```
[59]:      star_rating      title content_rating \
0          9.3      The Shawshank Redemption      R
1          9.2      The Godfather              R
2          9.1      The Godfather: Part II      R
4          8.9      Pulp Fiction                R
7          8.9  The Lord of the Rings: The Return of the King  PG-13

      genre  duration      actors_list
0    Crime    142  [u'Tim Robbins', u'Morgan Freeman', u'Bob Gunt...
1    Crime    175  [u'Marlon Brando', u'Al Pacino', u'James Caan']
2    Crime    200  [u'Al Pacino', u'Robert De Niro', u'Robert Duv...
4    Crime    154  [u'John Travolta', u'Uma Thurman', u'Samuel L...
7  Adventure    201  [u'Elijah Wood', u'Viggo Mortensen', u'Ian McK...
```

This works because tilde is the “not” operator in Python.

1.15 14. Filter a DataFrame by largest categories

Let’s say that you needed to filter the movies DataFrame by genre, but only include the 3 largest genres.

We’ll start by taking the `value_counts()` of genre and saving it as a Series called counts:

```
[60]: counts = movies.genre.value_counts()
counts
```

```
[60]: genre
Drama      278
Comedy     156
Action     136
Crime      124
Biography   77
Adventure   75
Animation   62
Horror      29
Mystery     16
Western      9
Sci-Fi       5
Thriller     5
Film-Noir    3
Family       2
History      1
Fantasy      1
Name: count, dtype: int64
```

The Series method `nlargest()` makes it easy to select the 3 largest values in this Series:

```
[61]: counts.nlargest(3)
```

```
[61]: genre
      Drama      278
      Comedy    156
      Action    136
      Name: count, dtype: int64
```

And all we actually need from this Series is the index:

```
[62]: counts.nlargest(3).index
```

```
[62]: Index(['Drama', 'Comedy', 'Action'], dtype='object', name='genre')
```

Finally, we can pass the index object to `isin()`, and it will be treated like a list of genres:

```
[63]: movies[movies.genre.isin(counts.nlargest(3).index)].head()
```

```
[63]:      star_rating      title \
3          9.0      The Dark Knight
5          8.9      12 Angry Men
9          8.9      Fight Club
11         8.8      Inception
12         8.8  Star Wars: Episode V - The Empire Strikes Back

      content_rating  genre  duration \
3          PG-13  Action      152
5      NOT RATED  Drama       96
9           R    Drama      139
11         PG-13  Action      148
12          PG   Action      124

      actors_list
3  [u'Christian Bale', u'Heath Ledger', u'Aaron E...
5  [u'Henry Fonda', u'Lee J. Cobb', u'Martin Bals...
9  [u'Brad Pitt', u'Edward Norton', u'Helena Bonh...
11 [u'Leonardo DiCaprio', u'Joseph Gordon-Levitt'...
12 [u'Mark Hamill', u'Harrison Ford', u'Carrie Fi...
```

Thus, only Drama and Comedy and Action movies remain in the DataFrame.

1.16 15. Handle missing values

Let's look at a dataset of UFO sightings:

```
[64]: ufo.head()
```

```
[64]:      City Colors Reported Shape Reported State \
0      Ithaca      NaN      TRIANGLE    NY
1  Willingboro      NaN      OTHER    NJ
2    Holyoke      NaN      OVAL    CO
```

3	Abilene	NaN	DISK	KS
4	New York Worlds Fair	NaN	LIGHT	NY

	Time
0	1930-06-01 22:00:00
1	1930-06-30 20:00:00
2	1931-02-15 14:00:00
3	1931-06-01 13:00:00
4	1933-04-18 19:00:00

You'll notice that some of the values are missing.

To find out how many values are missing in each column, you can use the `isna()` method and then take the `sum()`:

```
[65]: ufo.isna().sum()
```

```
[65]: City                26
Colors Reported    15359
Shape Reported    2644
State              0
Time              0
dtype: int64
```

`isna()` generated a DataFrame of True and False values, and `sum()` converted all of the True values to 1 and added them up.

Similarly, you can find out the percentage of values that are missing by taking the `mean()` of `isna()`:

```
[66]: ufo.isna().mean()
```

```
[66]: City                0.001425
Colors Reported    0.842004
Shape Reported    0.144948
State              0.000000
Time              0.000000
dtype: float64
```

If you want to drop the columns that have any missing values, you can use the `dropna()` method:

```
[67]: ufo.dropna(axis='columns').head()
```

```
[67]:   State                Time
0    NY 1930-06-01 22:00:00
1    NJ 1930-06-30 20:00:00
2    CO 1931-02-15 14:00:00
3    KS 1931-06-01 13:00:00
4    NY 1933-04-18 19:00:00
```

Or if you want to drop columns in which more than 10% of the values are missing, you can set a threshold for `dropna()`:

```
[68]: ufo.dropna(thresh=len(ufo)*0.9, axis='columns').head()
```

```
[68]:
```

	City	State		Time
0	Ithaca	NY	1930-06-01	22:00:00
1	Willingboro	NJ	1930-06-30	20:00:00
2	Holyoke	CO	1931-02-15	14:00:00
3	Abilene	KS	1931-06-01	13:00:00
4	New York Worlds Fair	NY	1933-04-18	19:00:00

`len(ufo)` returns the total number of rows, and then we multiply that by 0.9 to tell pandas to only keep columns in which at least 90% of the values are not missing.

1.17 16. Split a string into multiple columns

Let's create another example DataFrame:

```
[69]: df = pd.DataFrame({'name': ['John Arthur Doe', 'Jane Ann Smith'],  
                        'location': ['Los Angeles, CA', 'Washington, DC']})  
df
```

```
[69]:
```

	name	location
0	John Arthur Doe	Los Angeles, CA
1	Jane Ann Smith	Washington, DC

What if we wanted to split the “name” column into three separate columns, for first, middle, and last name? We would use the `str.split()` method and tell it to split on a space character and expand the results into a DataFrame:

```
[70]: df.name.str.split(' ', expand=True)
```

```
[70]:
```

	0	1	2
0	John	Arthur	Doe
1	Jane	Ann	Smith

These three columns can actually be saved to the original DataFrame in a single assignment statement:

```
[71]: df[['first', 'middle', 'last']] = df.name.str.split(' ', expand=True)  
df
```

```
[71]:
```

	name	location	first	middle	last
0	John Arthur Doe	Los Angeles, CA	John	Arthur	Doe
1	Jane Ann Smith	Washington, DC	Jane	Ann	Smith

What if we wanted to split a string, but only keep one of the resulting columns? For example, let's split the location column on “comma space”:

```
[72]: df.location.str.split(', ', expand=True)
```

```
[72]:
```

	0	1
0	Los Angeles	CA
1	Washington	DC

If we only cared about saving the city name in column 0, we can just select that column and save it to the DataFrame:

```
[73]: df['city'] = df.location.str.split(' ', expand=True)[0]
df
```

```
[73]:
```

	name	location	first	middle	last	city
0	John Arthur Doe	Los Angeles, CA	John	Arthur	Doe	Los Angeles
1	Jane Ann Smith	Washington, DC	Jane	Ann	Smith	Washington

1.18 17. Expand a Series of lists into a DataFrame

Let's create another example DataFrame:

```
[74]: df = pd.DataFrame({'col_one': ['a', 'b', 'c'], 'col_two': [[10, 40], [20, 50], [30, 60]]})
df
```

```
[74]:
```

	col_one	col_two
0	a	[10, 40]
1	b	[20, 50]
2	c	[30, 60]

There are two columns, and the second column contains regular Python lists of integers.

If we wanted to expand the second column into its own DataFrame, we can use the `apply()` method on that column and pass it the `Series()` constructor:

```
[75]: df_new = df.col_two.apply(pd.Series)
df_new
```

```
[75]:
```

	0	1
0	10	40
1	20	50
2	30	60

And by using the `concat()` function, you can combine the original DataFrame with the new DataFrame:

```
[76]: pd.concat([df, df_new], axis='columns')
```

```
[76]:
```

	col_one	col_two	0	1
0	a	[10, 40]	10	40
1	b	[20, 50]	20	50
2	c	[30, 60]	30	60

1.19 18. Aggregate by multiple functions

Let's look at a DataFrame of orders from the Chipotle restaurant chain:

```
[77]: orders.head(10)
```

```
[77]:
```

	order_id	quantity	item_name \
0	1	1	Chips and Fresh Tomato Salsa
1	1	1	Izze
2	1	1	Nantucket Nectar
3	1	1	Chips and Tomatillo-Green Chili Salsa
4	2	2	Chicken Bowl
5	3	1	Chicken Bowl
6	3	1	Side of Chips
7	4	1	Steak Burrito
8	4	1	Steak Soft Tacos
9	5	1	Steak Burrito

	choice_description	item_price
0	NaN	2.39
1	[Clementine]	3.39
2	[Apple]	3.39
3	NaN	2.39
4	[Tomatillo-Red Chili Salsa (Hot), [Black Beans...	16.98
5	[Fresh Tomato Salsa (Mild), [Rice, Cheese, Sou...	10.98
6	NaN	1.69
7	[Tomatillo Red Chili Salsa, [Fajita Vegetables...	11.75
8	[Tomatillo Green Chili Salsa, [Pinto Beans, Ch...	9.25
9	[Fresh Tomato Salsa, [Rice, Black Beans, Pinto...	9.25

Each order has an `order_id` and consists of one or more rows. To figure out the total price of an order, you sum the `item_price` for that `order_id`. For example, here's the total price of order number 1:

```
[78]: orders[orders.order_id == 1].item_price.sum()
```

```
[78]: 11.56
```

If you wanted to calculate the total price of every order, you would `groupby()` `order_id` and then take the sum of `item_price` for each group:

```
[79]: orders.groupby('order_id').item_price.sum().head()
```

```
[79]: order_id
1    11.56
2    16.98
3    12.67
4    21.00
5    13.70
```

Name: item_price, dtype: float64

However, you're not actually limited to aggregating by a single function such as `sum()`. To aggregate by multiple functions, you use the `agg()` method and pass it a list of functions such as `sum()` and `count()`:

```
[80]: orders.groupby('order_id').item_price.agg(['sum', 'count']).head()
```

```
[80]:
```

	sum	count
order_id		
1	11.56	4
2	16.98	1
3	12.67	2
4	21.00	2
5	13.70	2

That gives us the total price of each order as well as the number of items in each order.

1.20 19. Combine the output of an aggregation with a DataFrame

Let's take another look at the orders DataFrame:

```
[81]: orders.head(10)
```

```
[81]:
```

	order_id	quantity	item_name \
0	1	1	Chips and Fresh Tomato Salsa
1	1	1	Izze
2	1	1	Nantucket Nectar
3	1	1	Chips and Tomatillo-Green Chili Salsa
4	2	2	Chicken Bowl
5	3	1	Chicken Bowl
6	3	1	Side of Chips
7	4	1	Steak Burrito
8	4	1	Steak Soft Tacos
9	5	1	Steak Burrito

	choice_description	item_price
0	NaN	2.39
1	[Clementine]	3.39
2	[Apple]	3.39
3	NaN	2.39
4	[Tomatillo-Red Chili Salsa (Hot), [Black Beans...	16.98
5	[Fresh Tomato Salsa (Mild), [Rice, Cheese, Sou...	10.98
6	NaN	1.69
7	[Tomatillo Red Chili Salsa, [Fajita Vegetables...	11.75
8	[Tomatillo Green Chili Salsa, [Pinto Beans, Ch...	9.25
9	[Fresh Tomato Salsa, [Rice, Black Beans, Pinto...	9.25

What if we wanted to create a new column listing the total price of each order? Recall that we

calculated the total price using the `sum()` method:

```
[82]: orders.groupby('order_id').item_price.sum().head()
```

```
[82]: order_id
1      11.56
2      16.98
3      12.67
4      21.00
5      13.70
Name: item_price, dtype: float64
```

`sum()` is an aggregation function, which means that it returns a reduced version of the input data.

In other words, the output of the `sum()` function:

```
[83]: len(orders.groupby('order_id').item_price.sum())
```

```
[83]: 1834
```

...is smaller than the input to the function:

```
[84]: len(orders.item_price)
```

```
[84]: 4622
```

The solution is to use the `transform()` method, which performs the same calculation but returns output data that is the same shape as the input data:

```
[85]: total_price = orders.groupby('order_id').item_price.transform('sum')
len(total_price)
```

```
[85]: 4622
```

We'll store the results in a new DataFrame column called `total_price`:

```
[86]: orders['total_price'] = total_price
orders.head(10)
```

```
[86]:
```

	order_id	quantity	item_name \
0	1	1	Chips and Fresh Tomato Salsa
1	1	1	Izze
2	1	1	Nantucket Nectar
3	1	1	Chips and Tomatillo-Green Chili Salsa
4	2	2	Chicken Bowl
5	3	1	Chicken Bowl
6	3	1	Side of Chips
7	4	1	Steak Burrito
8	4	1	Steak Soft Tacos
9	5	1	Steak Burrito

	choice_description	item_price	total_price
0	NaN	2.39	11.56
1	[Clementine]	3.39	11.56
2	[Apple]	3.39	11.56
3	NaN	2.39	11.56
4	[Tomatillo-Red Chili Salsa (Hot), [Black Beans...	16.98	16.98
5	[Fresh Tomato Salsa (Mild), [Rice, Cheese, Sou...	10.98	12.67
6	NaN	1.69	12.67
7	[Tomatillo Red Chili Salsa, [Fajita Vegetables...	11.75	21.00
8	[Tomatillo Green Chili Salsa, [Pinto Beans, Ch...	9.25	21.00
9	[Fresh Tomato Salsa, [Rice, Black Beans, Pinto...	9.25	13.70

As you can see, the total price of each order is now listed on every single line.

That makes it easy to calculate the percentage of the total order price that each line represents:

```
[87]: orders['percent_of_total'] = orders.item_price / orders.total_price
orders.head(10)
```

```
[87]:  order_id  quantity  item_name \
0         1         1  Chips and Fresh Tomato Salsa
1         1         1                Izze
2         1         1  Nantucket Nectar
3         1         1  Chips and Tomatillo-Green Chili Salsa
4         2         2  Chicken Bowl
5         3         1  Chicken Bowl
6         3         1  Side of Chips
7         4         1  Steak Burrito
8         4         1  Steak Soft Tacos
9         5         1  Steak Burrito
```

	choice_description	item_price	total_price
0	NaN	2.39	11.56
1	[Clementine]	3.39	11.56
2	[Apple]	3.39	11.56
3	NaN	2.39	11.56
4	[Tomatillo-Red Chili Salsa (Hot), [Black Beans...	16.98	16.98
5	[Fresh Tomato Salsa (Mild), [Rice, Cheese, Sou...	10.98	12.67
6	NaN	1.69	12.67
7	[Tomatillo Red Chili Salsa, [Fajita Vegetables...	11.75	21.00
8	[Tomatillo Green Chili Salsa, [Pinto Beans, Ch...	9.25	21.00
9	[Fresh Tomato Salsa, [Rice, Black Beans, Pinto...	9.25	13.70

	percent_of_total
0	0.206747
1	0.293253
2	0.293253
3	0.206747

```

4      1.000000
5      0.866614
6      0.133386
7      0.559524
8      0.440476
9      0.675182

```

1.21 20. Select a slice of rows and columns

Let's take a look at another dataset:

```
[88]: titanic.head()
```

```
[88]:
```

	PassengerId	Survived	Pclass	\
0	1	0	3	
1	2	1	1	
2	3	1	3	
3	4	1	1	
4	5	0	3	

	Name	Sex	Age	SibSp	\
0	Braund, Mr. Owen Harris	male	22.0	1	
1	Cumings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	
2	Heikkinen, Miss. Laina	female	26.0	0	
3	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	
4	Allen, Mr. William Henry	male	35.0	0	

	Parch	Ticket	Fare	Cabin	Embarked
0	0	A/5 21171	7.2500	NaN	S
1	0	PC 17599	71.2833	C85	C
2	0	STON/O2. 3101282	7.9250	NaN	S
3	0	113803	53.1000	C123	S
4	0	373450	8.0500	NaN	S

This is the famous Titanic dataset, which shows information about passengers on the Titanic and whether or not they survived.

If you wanted a numerical summary of the dataset, you would use the `describe()` method:

```
[89]: titanic.describe()
```

```
[89]:
```

	PassengerId	Survived	Pclass	Age	SibSp	\
count	891.000000	891.000000	891.000000	714.000000	891.000000	
mean	446.000000	0.383838	2.308642	29.699118	0.523008	
std	257.353842	0.486592	0.836071	14.526497	1.102743	
min	1.000000	0.000000	1.000000	0.420000	0.000000	
25%	223.500000	0.000000	2.000000	20.125000	0.000000	
50%	446.000000	0.000000	3.000000	28.000000	0.000000	
75%	668.500000	1.000000	3.000000	38.000000	1.000000	

```
max      891.000000      1.000000      3.000000      80.000000      8.000000
```

```

      Parch      Fare
count  891.000000  891.000000
mean    0.381594   32.204208
std     0.806057   49.693429
min     0.000000    0.000000
25%     0.000000    7.910400
50%     0.000000   14.454200
75%     0.000000   31.000000
max     6.000000  512.329200

```

However, the resulting DataFrame might be displaying more information than you need.

If you wanted to filter it to only show the “five-number summary”, you can use the `loc` accessor and pass it a slice of the “min” through the “max” row labels:

```
[90]: titanic.describe().loc['min':'max']
```

```

[90]:      PassengerId  Survived  Pclass     Age  SibSp  Parch      Fare
min           1.0         0.0      1.0   0.420    0.0    0.0     0.0000
25%          223.5         0.0      2.0  20.125    0.0    0.0     7.9104
50%          446.0         0.0      3.0  28.000    0.0    0.0    14.4542
75%          668.5         1.0      3.0  38.000    1.0    0.0    31.0000
max          891.0         1.0      3.0  80.000    8.0    6.0   512.3292

```

And if you’re not interested in all of the columns, you can also pass it a slice of column labels:

```
[91]: titanic.describe().loc['min':'max', 'Pclass':'Parch']
```

```

[91]:      Pclass     Age  SibSp  Parch
min       1.0   0.420    0.0    0.0
25%       2.0  20.125    0.0    0.0
50%       3.0  28.000    0.0    0.0
75%       3.0  38.000    1.0    0.0
max       3.0  80.000    8.0    6.0

```

1.22 21. Reshape a MultiIndexed Series

The Titanic dataset has a “Survived” column made up of ones and zeros, so you can calculate the overall survival rate by taking a `mean()` of that column:

```
[92]: titanic.Survived.mean()
```

```
[92]: 0.3838383838383838
```

If you wanted to calculate the survival rate by a single category such as “Sex”, you would use a `groupby()`:

```
[93]: titanic.groupby('Sex').Survived.mean()
```

```
[93]: Sex
      female    0.742038
      male      0.188908
      Name: Survived, dtype: float64
```

And if you wanted to calculate the survival rate across two different categories at once, you would `groupby()` both of those categories:

```
[94]: titanic.groupby(['Sex', 'Pclass']).Survived.mean()
```

```
[94]: Sex      Pclass
      female  1          0.968085
              2          0.921053
              3          0.500000
      male    1          0.368852
              2          0.157407
              3          0.135447
      Name: Survived, dtype: float64
```

This shows the survival rate for every combination of Sex and Passenger Class. It's stored as a MultiIndexed Series, meaning that it has multiple index levels to the left of the actual data.

It can be hard to read and interact with data in this format, so it's often more convenient to reshape a MultiIndexed Series into a DataFrame by using the `unstack()` method:

```
[95]: titanic.groupby(['Sex', 'Pclass']).Survived.mean().unstack()
```

```
[95]: Pclass      1          2          3
      Sex
      female  0.968085  0.921053  0.500000
      male    0.368852  0.157407  0.135447
```

This DataFrame contains the same exact data as the MultiIndexed Series, except that now you can interact with it using familiar DataFrame methods.

1.23 22. Create a pivot table

If you often create DataFrames like the one above, you might find it more convenient to use the `pivot_table()` method instead:

```
[96]: titanic.pivot_table(index='Sex', columns='Pclass', values='Survived',
      ↪aggfunc='mean')
```

```
[96]: Pclass      1          2          3
      Sex
      female  0.968085  0.921053  0.500000
      male    0.368852  0.157407  0.135447
```

With a pivot table, you directly specify the index, the columns, the values, and the aggregation function.

An added benefit of a pivot table is that you can easily add row and column totals by setting `margins=True`:

```
[97]: titanic.pivot_table(index='Sex', columns='Pclass', values='Survived',  
    ↪aggfunc='mean',  
    margins=True)
```

```
[97]: Pclass      1      2      3      All  
Sex  
female  0.968085  0.921053  0.500000  0.742038  
male    0.368852  0.157407  0.135447  0.188908  
All     0.629630  0.472826  0.242363  0.383838
```

This shows the overall survival rate as well as the survival rate by Sex and Passenger Class.

Finally, you can create a cross-tabulation just by changing the aggregation function from “mean” to “count”:

```
[98]: titanic.pivot_table(index='Sex', columns='Pclass', values='Survived',  
    ↪aggfunc='count',  
    margins=True)
```

```
[98]: Pclass    1    2    3  All  
Sex  
female   94   76  144  314  
male    122  108  347  577  
All     216  184  491  891
```

This shows the number of records that appear in each combination of categories.

1.24 23. Convert continuous data into categorical data

Let's take a look at the Age column from the Titanic dataset:

```
[99]: titanic.Age.head(10)
```

```
[99]: 0    22.0  
1    38.0  
2    26.0  
3    35.0  
4    35.0  
5     NaN  
6    54.0  
7     2.0  
8    27.0  
9    14.0  
Name: Age, dtype: float64
```

It's currently continuous data, but what if you wanted to convert it into categorical data?

One solution would be to label the age ranges, such as “child”, “young adult”, and “adult”. The best way to do this is by using the `cut()` function:

```
[100]: pd.cut(titanic.Age, bins=[0, 18, 25, 99], labels=['child', 'young adult', 'adult']).head(10)
```

```
[100]: 0    young adult
      1         adult
      2         adult
      3         adult
      4         adult
      5          NaN
      6         adult
      7         child
      8         adult
      9         child
      Name: Age, dtype: category
      Categories (3, object): ['child' < 'young adult' < 'adult']
```

This assigned each value to a bin with a label. Ages 0 to 18 were assigned the label “child”, ages 18 to 25 were assigned the label “young adult”, and ages 25 to 99 were assigned the label “adult”.

Notice that the data type is now “category”, and the categories are automatically ordered.

1.25 24. Change display options

Let’s take another look at the Titanic dataset:

```
[101]: titanic.head()
```

```
[101]: PassengerId  Survived  Pclass  \
0             1         0         3
1             2         1         1
2             3         1         3
3             4         1         1
4             5         0         3

      Name      Sex  Age  SibSp  \
0  Braund, Mr. Owen Harris    male  22.0      1
1  Cumings, Mrs. John Bradley (Florence Briggs Th...  female  38.0      1
2    Heikkinen, Miss. Laina   female  26.0      0
3  Futrelle, Mrs. Jacques Heath (Lily May Peel)   female  35.0      1
4    Allen, Mr. William Henry    male  35.0      0

      Parch      Ticket    Fare Cabin Embarked
0         0   A/5 21171    7.2500   NaN        S
1         0   PC 17599   71.2833   C85        C
2         0  STON/O2. 3101282   7.9250   NaN        S
3         0   113803   53.1000  C123        S
```

```
4      0      373450    8.0500    NaN      S
```

Notice that the Age column has 1 decimal place and the Fare column has 4 decimal places. What if you wanted to standardize the display to use 2 decimal places?

You can use the `set_option()` function:

```
[102]: pd.set_option('display.float_format', '{:.2f}'.format)
```

The first argument is the name of the option, and the second argument is a Python format string.

```
[103]: titanic.head()
```

```
[103]:
```

	PassengerId	Survived	Pclass	\
0	1	0	3	
1	2	1	1	
2	3	1	3	
3	4	1	1	
4	5	0	3	

	Name	Sex	Age	SibSp	\
0	Braund, Mr. Owen Harris	male	22.00	1	
1	Cumings, Mrs. John Bradley (Florence Briggs Th...	female	38.00	1	
2	Heikkinen, Miss. Laina	female	26.00	0	
3	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.00	1	
4	Allen, Mr. William Henry	male	35.00	0	

	Parch	Ticket	Fare	Cabin	Embarked
0	0	A/5 21171	7.25	NaN	S
1	0	PC 17599	71.28	C85	C
2	0	STON/O2. 3101282	7.92	NaN	S
3	0	113803	53.10	C123	S
4	0	373450	8.05	NaN	S

You can see that Age and Fare are now using 2 decimal places. Note that this did not change the underlying data, only the display of the data.

You can also reset any option back to its default:

```
[104]: pd.reset_option('display.float_format')
```

There are many more options you can specify in a similar way.

1.26 25. Style a DataFrame

The previous trick is useful if you want to change the display of your entire notebook. However, a more flexible and powerful approach is to define the style of a particular DataFrame.

Let's return to the stocks DataFrame:

```
[105]: stocks
```



```
[105]:
```

	Date	Close	Volume	Symbol
0	2016-10-03	31.50	14070500	CSCO
1	2016-10-03	112.52	21701800	AAPL
2	2016-10-03	57.42	19189500	MSFT
3	2016-10-04	113.00	29736800	AAPL
4	2016-10-04	57.24	20085900	MSFT
5	2016-10-04	31.35	18460400	CSCO
6	2016-10-05	57.64	16726400	MSFT
7	2016-10-05	31.59	11808600	CSCO
8	2016-10-05	113.05	21453100	AAPL

We can create a dictionary of format strings that specifies how each column should be formatted:

```
[106]: format_dict = {'Date': '{:%m/%d/%y}', 'Close': '${:.2f}', 'Volume': '{:,}'}
```

And then we can pass it to the DataFrame's `style.format()` method:

```
[107]: stocks.style.format(format_dict)
```

```
[107]: <pandas.io.formats.style.Styler at 0x13b170490>
```

Notice that the Date is now in month-day-year format, the closing price has a dollar sign, and the Volume has commas.

We can apply more styling by chaining additional methods:

```
[108]: (stocks.style.format(format_dict)
        .hide(axis='index')
        .highlight_min('Close', color='red')
        .highlight_max('Close', color='lightgreen')
        )
```

```
[108]: <pandas.io.formats.style.Styler at 0x13c6b03d0>
```

We've now hidden the index, highlighted the minimum Close value in red, and highlighted the maximum Close value in green.

Here's another example of DataFrame styling:

```
[109]: (stocks.style.format(format_dict)
        .hide(axis='index')
        .background_gradient(subset='Volume', cmap='Blues')
        )
```

```
[109]: <pandas.io.formats.style.Styler at 0x13b1c8e50>
```

The Volume column now has a background gradient to help you easily identify high and low values.

And here's one final example:

```
[110]: (stocks.style.format(format_dict)
        .hide(axis='index')
        .bar('Volume', color='lightblue', align='zero')
        .set_caption('Stock Prices from October 2016')
        )
```

```
[110]: <pandas.io.formats.style.Styler at 0x13c707250>
```

There's now a bar chart within the Volume column and a caption above the DataFrame.

Note that there are many more options for how you can style your DataFrame.

1.27 Bonus: Profile a DataFrame

Let's say that you've got a new dataset, and you want to quickly explore it without too much work. There's a separate package called [ydata-profiling](#) that is designed for this purpose.

First you have to install it using conda or pip. Once that's done, you import `ydata_profiling`:

```
[111]: import ydata_profiling
```

Then, simply run the `ProfileReport()` function and pass it any DataFrame. It returns an interactive HTML report:

- The first section is an overview of the dataset and a list of possible issues with the data.
- The next section gives a summary of each column. You can click “toggle details” for even more information.
- The third section shows a heatmap of the correlation between columns.
- And the fourth section shows the head of the dataset.

```
[112]: ydata_profiling.ProfileReport(titanic)
```

```
Summarize dataset:  0%|          | 0/5 [00:00<?, ?it/s]
```

```
Generate report structure:  0%|          | 0/1 [00:00<?, ?it/s]
```

```
Render HTML:  0%|          | 0/1 [00:00<?, ?it/s]
```

```
<IPython.core.display.HTML object>
```

```
[112]:
```