

Introduction to Pandas in Python

Estimated time needed: **15** minutes

Objectives

After completing this lab you will be able to:

- Use Pandas to access and view data

The table has one row for each album and several columns.

You can see the dataset here:

Artist	Album	Released	Length	Genre	Music recording sales (millions)	Claimed sales (millions)
Michael Jackson	Thriller	1982	00:42:19	Pop, rock, R&B	46	65
AC/DC	Back in Black	1980	00:42:11	Hard rock	26.1	50
Pink Floyd	The Dark Side of the Moon	1973	00:42:49	Progressive rock	24.2	45
Whitney Houston	The Bodyguard	1992	00:57:44	Soundtrack/R&B, soul, pop	26.1	50
Meat Loaf	Bat Out of Hell	1977	00:46:33	Hard rock, progressive rock	20.6	43
Eagles	Their Greatest Hits (1971-1975)	1976	00:43:08	Rock, soft rock, folk rock	32.2	42
Bee Gees	Saturday Night Fever	1977	1:15:54	Disco	20.6	40
Fleetwood Mac	Rumours	1977	00:40:01	Soft rock	27.9	40

```
# Dependency needed to install file
```

```
!pip install xlrd
!pip install openpyxl
```

```
Requirement already satisfied: xlrd in
/home/jupyterlab/conda/envs/python/lib/python3.7/site-packages (1.2.0)
Requirement already satisfied: openpyxl in
/home/jupyterlab/conda/envs/python/lib/python3.7/site-packages (3.1.2)
Requirement already satisfied: et-xmlfile in
/home/jupyterlab/conda/envs/python/lib/python3.7/site-packages (from
openpyxl) (1.1.0)
```

```
# Import required library
```

```
import pandas as pd
```

After the import command, we now have access to a large number of pre-built classes and functions. This assumes the library is installed; in our lab environment all the necessary libraries are installed. One way pandas allow you to work with data is a dataframe. Let's go through the process to go from a comma separated values (.csv) file to a dataframe. This variable `csv_path` stores the path of the .csv, which is used as an argument to the `read_csv` function. The result is stored in the object `df`, this is a common short form used for a variable referring to a Pandas dataframe.

```
# Read data from CSV file
```

```
csv_path = 'https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBMDeveloperSkillsNetwork-PY0101EN-SkillsNetwork/labs/Module%204/data/TopSellingAlbums.csv'
df = pd.read_csv(csv_path)
```

We can use the method head() to examine the first five rows of a dataframe:

```
# Print first five rows of the dataframe
```

```
df.head()
```

	Artist	Album	Released	Length	\
0	Michael Jackson	Thriller	1982	0:42:19	
1	AC/DC	Back in Black	1980	0:42:11	
2	Pink Floyd	The Dark Side of the Moon	1973	0:42:49	
3	Whitney Houston	The Bodyguard	1992	0:57:44	
4	Meat Loaf	Bat Out of Hell	1977	0:46:33	

	Genre	Music Recording Sales (millions)	\
0	pop, rock, R&B	46.0	
1	hard rock	26.1	
2	progressive rock	24.2	
3	R&B, soul, pop	27.4	
4	hard rock, progressive rock	20.6	

	Claimed Sales (millions)	Released.1	Soundtrack	Rating
0	65	30-Nov-82	NaN	10.0
1	50	25-Jul-80	NaN	9.5
2	45	01-Mar-73	NaN	9.0
3	44	17-Nov-92	Y	8.5
4	43	21-Oct-77	NaN	8.0

We use the path of the excel file and the function read_excel. The result is a data frame as before:

```
# Read data from Excel File and print the first five rows
```

```
xlsx_path = 'https://s3-api.us-geo.objectstorage.softlayer.net/cf-courses-data/CognitiveClass/PY0101EN/Chapter%204/Datasets/TopSellingAlbums.xlsx'
```

```
df = pd.read_excel(xlsx_path)
df.head()
```

	Artist	Album	Released	Length	\
0	Michael Jackson	Thriller	1982	00:42:19	
1	AC/DC	Back in Black	1980	00:42:11	
2	Pink Floyd	The Dark Side of the Moon	1973	00:42:49	

3	Whitney Houston	The Bodyguard	1992	00:57:44
4	Meat Loaf	Bat Out of Hell	1977	00:46:33

	Genre	Music Recording Sales (millions)	\
0	pop, rock, R&B	46.0	
1	hard rock	26.1	
2	progressive rock	24.2	
3	R&B, soul, pop	27.4	
4	hard rock, progressive rock	20.6	

	Claimed Sales (millions)	Released.1	Soundtrack	Rating
0	65	1982-11-30	NaN	10.0
1	50	1980-07-25	NaN	9.5
2	45	1973-03-01	NaN	9.0
3	44	1992-11-17	Y	8.5
4	43	1977-10-21	NaN	8.0

We can access the column Length and assign it a new dataframe x:

```
# Access to the column Length
```

```
x = df[['Length']]
x
```

```

      Length
0  00:42:19
1  00:42:11
2  00:42:49
3  00:57:44
4  00:46:33
5  00:43:08
6  01:15:54
7  00:40:01
```

The process is shown in the figure:

You can also get a column as a series. You can think of a Pandas series as a 1-D dataframe. Just use one bracket:

```
# Get the column as a series
```

```
x = df['Length']
x
```

```

0    00:42:19
1    00:42:11
2    00:42:49
3    00:57:44
```

```
4    00:46:33
5    00:43:08
6    01:15:54
7    00:40:01
Name: Length, dtype: object
```

You can also get a column as a dataframe. For example, we can assign the column Artist:

```
# Get the column as a dataframe
```

```
x = df[['Artist']]
type(x)
```

```
pandas.core.frame.DataFrame
```

You can do the same thing for multiple columns; we just put the dataframe name, in this case, df, and the name of the multiple column headers enclosed in double brackets. The result is a new dataframe comprised of the specified columns:

```
# Access to multiple columns
```

```
y = df[['Artist', 'Length', 'Genre']]
y
```

	Artist	Length	Genre
0	Michael Jackson	00:42:19	pop, rock, R&B
1	AC/DC	00:42:11	hard rock
2	Pink Floyd	00:42:49	progressive rock
3	Whitney Houston	00:57:44	R&B, soul, pop
4	Meat Loaf	00:46:33	hard rock, progressive rock
5	Eagles	00:43:08	rock, soft rock, folk rock
6	Bee Gees	01:15:54	disco
7	Fleetwood Mac	00:40:01	soft rock

The process is shown in the figure:

One way to access unique elements is the iloc method, where you can access the 1st row and the 1st column as follows:

```
# Access the value on the first row and the first column
```

```
df.iloc[0, 0]
```

```
'Michael Jackson'
```

You can access the 2nd row and the 1st column as follows:

```
# Access the value on the second row and the first column  
df.iloc[1,0]  
'AC/DC'
```

You can access the 1st row and the 3rd column as follows:

```
# Access the value on the first row and the third column  
df.iloc[0,2]  
1982  
# Access the value on the second row and the third column  
df.iloc[1,2]  
1980
```

This is shown in the following image

You can access the column using the name as well, the following are the same as above:

```
# Access the column using the name  
df.loc[0, 'Artist']  
'Michael Jackson'  
# Access the column using the name  
df.loc[1, 'Artist']  
'AC/DC'  
# Access the column using the name  
df.loc[0, 'Released']  
1982  
# Access the column using the name  
df.loc[1, 'Released']  
1980
```

You can perform slicing using both the index and the name of the column:

```
# Slicing the dataframe
```

```
df.iloc[0:2, 0:3]
```

	Artist	Album	Released
0	Michael Jackson	Thriller	1982
1	AC/DC	Back in Black	1980

```
# Slicing the dataframe using name
```

```
df.loc[0:2, 'Artist':'Released']
```

	Artist	Album	Released
0	Michael Jackson	Thriller	1982
1	AC/DC	Back in Black	1980
2	Pink Floyd	The Dark Side of the Moon	1973

Use a variable q to store the column Rating as a dataframe

```
# Write your code below and press Shift+Enter to execute
```

```
q = df[['Rating']]
```

```
q
```

	Rating
0	10.0
1	9.5
2	9.0
3	8.5
4	8.0
5	7.5
6	7.0
7	6.5

Assign the variable q to the dataframe that is made up of the column Released and Artist:

```
# Write your code below and press Shift+Enter to execute
```

```
q = df[['Released', 'Artist', 'Rating']]
```

```
q
```

	Released	Artist	Rating
0	1982	Michael Jackson	10.0
1	1980	AC/DC	9.5
2	1973	Pink Floyd	9.0
3	1992	Whitney Houston	8.5

4	1977	Meat Loaf	8.0
5	1976	Eagles	7.5
6	1977	Bee Gees	7.0
7	1977	Fleetwood Mac	6.5

Access the 2nd row and the 3rd column of df:

Write your code below and press Shift+Enter to execute

Use the following list to convert the dataframe index df to characters and assign it to df_new; find the element corresponding to the row index a and column 'Artist'. Then select the rows a through d for the column 'Artist'

```
new_index=['a','b','c','d','e','f','g','h']
```

Congratulations, you have completed your first lesson and hands-on lab in Python.

Authors:

[Joseph Santarcangelo](#)

Joseph Santarcangelo has a PhD in Electrical Engineering, his research focused on using machine learning, signal processing, and computer vision to determine how videos impact human cognition. Joseph has been working for IBM since he completed his PhD.

Change Log

Date (YYYY-MM-DD)	Version	Changed By	Change Description
2022-01-10	2.1	Malika	Removed the readme for GitShare
2020-08-26	2.0	Lavanya	Moved lab to course repo in GitLab
2020-11-24	3.0	Nayef	Added new images

© IBM Corporation 2020. All rights reserved.